

<b>Manuscript Number:</b>	GIGA-D-20-00244R2	
<b>Full Title:</b>	Torix Rickettsia are widespread in arthropods and reflect a neglected symbiosis	
<b>Article Type:</b>	Research	
<b>Funding Information:</b>	Biotechnology and Biological Sciences Research Council (BB/M011186/1)	Dr. Jack Pilgrim
	Natural Environment Research Council (NE/L002450/1)	Ms. Helen R. Davison
<b>Abstract:</b>	<p><b>Background</b></p> <p>Rickettsia are intracellular bacteria best known as the causative agents of human and animal diseases. Although these medically important Rickettsia are often transmitted via haematophagous arthropods, other Rickettsia, such as those in the Torix group, appear to reside exclusively in invertebrates and protists with no secondary vertebrate host. Importantly, little is known about the diversity or host range of Torix group Rickettsia.</p> <p><b>Results</b></p> <p>This study describes the serendipitous discovery of Rickettsia amplicons in the Barcode of Life Data System (BOLD), a sequence database specifically designed for the curation of mtDNA barcodes. Out of 184,585 barcode sequences analysed, Rickettsia is observed in approximately 0.41% of barcode submissions and is more likely to be found than Wolbachia (0.17%). The Torix group of Rickettsia are shown to account for 95% of all unintended amplifications from the genus. A further targeted PCR screen of 1,612 individuals from 169 terrestrial and aquatic invertebrate species identified mostly Torix strains and supports the 'aquatic hot spot' hypothesis for Torix infection. Furthermore, the analysis of 1,341 Sequence Read Archive (SRA) deposits indicates Torix infections represent a significant proportion of all Rickettsia symbioses found in arthropod genome projects.</p> <p><b>Conclusions</b></p> <p>This study supports a previous hypothesis which suggests Torix Rickettsia are overrepresented in aquatic insects. In addition, multiple methods reveal further putative hot spots of Torix Rickettsia infection; including in phloem-feeding bugs, parasitoid wasps, spiders, and vectors of disease. The unknown host effects and transmission strategies of these endosymbionts make these newly discovered associations important to inform future directions of investigation involving the understudied Torix Rickettsia.</p>	
<b>Corresponding Author:</b>	Jack Pilgrim University of Liverpool Institute of Infection Veterinary and Ecological Sciences Neston, Cheshire, Liverpool UNITED KINGDOM	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	University of Liverpool Institute of Infection Veterinary and Ecological Sciences	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Jack Pilgrim	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Jack Pilgrim	
	Panupong Thongprem	
	Helen R. Davison	

	Stefanos Siozios
	Matthew Baylis
	Evgeny V. Zakharov
	Sujeevan Ratnasingham
	Jeremy R. deWaard
	Craig R. Macadam
	M. Alex Smith
	Gregory D. D. Hurst
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	<p>We would like to thank the reviewer for their comments. Please find below a point-by-point response from the authors.</p> <p>“The alignments seem to interchange N/n/-, but I checked the IQTree manual and IQTree treats them all the same. But I don't see how IQ tree handles phylogenetic inference at positions with N/n/-, which it calls "unknowns". For example, in the multigene alignment, I don't believe there is any position without gaps, so presumably it handles them. I'm assuming the alignments I'm looking at are the ones fed into IQTree and not the ones coming out of MAFFT and before Gblocks, but it might be good for you to confirm using this example: looking at the alignment BOLD_just_Rciektsia_COI_contaminants_alignment.fas, these two sequences have absolutely no overlap, and are significant truncated relative to the complete alignment (and these aren't the only two with this problem). How is the algorithm dealing with this, and does it introduce any artifacts? Do you get the same result if you remove sequences like these and use only sequences where all positions of the alignment have a character (both by removing short sequences like this and also trimming the ends of the alignment)?”</p> <p>The reasoning behind including sequences with missing character data in our alignments is based on previous work demonstrating that missing data in most cases should not decrease phylogenetic resolution for taxa with complete data (Wiens 2006, DOI: 10.1016/j.jbi.2005.04.001). To confirm this, we reran a modified alignment of the BOLD_just_Rickettsia_COI_contaminants_alignment.fas file by trimming ends by 50 nucleotides and removing any remaining truncated sequences to get rid of missing data (169 of the original 807 sequences removed), as suggested by the reviewer. In accordance with Wiens' observations, the generated tree (see FTP file 'BOLD Rickettsia trimmed.png') placed taxa into the same designated groups as the phylogeny with missing data in Figure 3 (Designations can be found at DOI:10.6084/m9.figshare.12801107). Additionally, the study cited above ran simulations to show that highly incomplete data can be accurately placed in phylogenetic trees as long as at least 50% of sequences contain complete data. Furthermore, it is suggested that the inclusion of sequences with missing data can also sometimes be better than exclusion as this additional data can subdivide misleading long branches. We have now included the reasoning behind using incomplete data in the methods section (lines 480-483).</p> <p>“How do double infections confound the results? Do they behave erratically in the cladograms, like a chimera would? Can this be a wider problem in interpreting the results? It seems like it would have to be, but you have chimeras you can use to examine this. Using those, is there a way to address this? With SRA data it seems like you should be able to look for sequence heterogeneity in the reads. Not sure about BOLD data.”</p> <p>Where double peaks were observed in 10/753 Rickettsia-associated taxa from BOLD, the base call was designated as 'N' (See FTP file 'BOLD_multigene_Rickettsia_alignment.fas'). This prevents erroneous placement of chimeric strains on the phylogeny. For BOLD data we unfortunately cannot reconstruct trees by teasing apart the individual strains of the double-infections because we cannot know what phase the double-peaks are in.</p>

The use of 'N' characters at double-peak sites could lead to potential problems in the interpretation of these 10 taxa at terminal branches of the phylogeny but the placement as *Torix Rickettsia* is not likely to be affected. Furthermore, these double-infections are a minority of the total taxa meaning their effects on interpreting results are likely to be minimal.

“Line 254-261: I think you need to add in a correction for multiple testing. You do at least two tests that are in the main text, but it sounds from the response to reviewer's comments that you did many more than that and are only reporting the ones that are  $P < 0.05$ . However, you should report how many tests you did to find those two and adjust for multiple testing. Otherwise, if you do 20 tests, you would expect to have 1 that is "significant" (for more information on multiple testing: <http://www.biostathandbook.com/multiplecomparisons.html>). In addition, I think it is important for people to know what comparisons were done that were not significant as these are also results. Addressing multiple testing seems like an issue throughout. In the methods other statistical tests were clearly undertaken where there was multiple testing.”

Two Fisher's exact tests (aquatic vs terrestrial insects-1 controlled for insect order and 1 uncontrolled) were detailed in the main text and additional file 7, as these were the only taxonomically 'matched' pairs. However, one additional test was performed initially to compare terrestrial vs aquatic invertebrates in general which did not give a significant p-value due to a hotspot of *Rickettsia* in spiders, which are known to be a hotspot for all inherited symbionts tested to date (*Wolbachia*, *Spiroplasma*, *Rickettsia*: Goodacre et al. 2006, doi: 10.1111/j.1365-294X.2005.02802.x.; *Cardinium*: Duron et al. 2008, doi: 10.1111/j.1365-294X.2008.03689.x.). This detail has now been added in Additional file 7 and lines 266-269. Overall, only 3 tests were done (2 significant and 1 not significant) and this indicates that *Torix Rickettsia* are over-represented in aquatic insects but this may not be the case for invertebrates in general.

“Line 354: Several, maybe many, of the *Wolbachia* integrations have no mutations or frameshifts, particularly in insects. Those with frameshifts and mutations are easier to find and identify as integrations such that the number of integrations without frameshifts and mutations is likely an underestimate, particularly given how many groups are still screening *Wolbachia* sequences out before assembling insect genomes. I have no idea how often that happens with *Rickettsia*, but it seems like, particularly as more groups use tools like blobplots.”

We thank the reviewer for raising this issue. We have now put a caveat at the end of this sentence to indicate that despite no frameshifts or mutations, it is still possible the sequences from this study are host integrations (lines 376-378).

The problem is likely to be less for *Rickettsia* than *Wolbachia*, due to differences in the mode of vertical transmission. *Wolbachia* is present in the germline stem cell niche, such DNA from the symbiont is available for incorporation into the germline. *Rickettsia*, in contrast, usually invades the egg after meiosis, through the follicular epithelium. Thus, *Rickettsia* DNA is much less present in the germline of insects, making integration less likely.

“Line 366-379: This section still has issues with respect to the study design being secondary data analysis. These lines are in the discussion, it is the time to say things like on line 377 that the over-representation here in BOLD data (if that is the data you are referring to, because I can't remember which one was 17/19 and there is not here that clarifies) could be the result of an amplification bias—in not producing the host copy of the gene, amplifying the *Rickettsia* gene, or both. Those issues are profound in secondary data usage and need to be addressed head-on so that others who read the paper do not misconstrue the results. Likewise, the SRA data is not random, so I am not sure the statement on line 379-381 is correct, and at very least it needs qualifications. If it is correct, you need to better argue in the manuscript why it is correct, like that you used a sampling scheme to reduce bias, or something like that. Personally, I think it is better to acknowledge the limitations that try to justify, as even if you have a sampling scheme, it can be biased. The PCR screen listed in the table of this manuscripts seems biased from my

quick look (e.g. an over-representation of mosquitoes).”

The “17/19 strains” being Torix is a reference to the targeted screen (not the BOLD screen) which was used alongside the BOLD data because of the aforementioned biases relating to amplification bias and this has now been clarified on lines 401-402. Additionally, we have added a sentence to the results section explicitly quoting the 17 strains of Rickettsia found in the targeted screen (lines 248-250). Although 95% of Rickettsia amplifications from BOLD are Torix, we already mention that this is likely due to primer bias (lines 321-324). Subsequently, the targeted screen is used in part to negate the problems of relying entirely on secondary data.

Of course, many studies which aim to investigate the distribution of a symbiont will have sampling and methodological biases. However, having multiple screening strategies, as we have here, is likely to give a more nuanced and holistic view of Torix Rickettsia ecology. We believe that the combined use of several screening methods is a strength and not a weakness of the study. Despite this, we have now added a separate section detailing the limitations of the study (lines 358-388).

Specifically, regarding lines 379-381 (of the 1st revision), this statement is based not just on SRA data but also the targeted screen from this study and Weinert’s study as mentioned in the previous lines. Thus, the SRA is corroborating two separate targeted screens (one which lacked spiders and aquatic insects demonstrating a high number of Belli infections, and another which included spiders and aquatic insects demonstrating a high number of Torix infections.). Subsequently, for clarity we have now changed the statement “Our additional use of a bioinformatics approach based on the SRA appears to confirm that Belli and Torix are two of the most common Rickettsia groups among arthropods.” to “Our additional use of a bioinformatics approach based on the SRA appears to corroborate targeted screen data indicating that Belli and Torix are two of the most common Rickettsia groups among arthropods.” (lines 403-406).

“Line 387-388: please provide more details. I don’t remember reading that. Pointing to an exact result, for instance of how many strains of the same MLST type are in different insect orders is necessary. It should have been in the results if it is in the discussion. In fact, if I look at the figures, the Wolbachia in Figure 2 actually seem to be grouped by insect host taxa at this level. The same is true for Figure 3 for the Rickettsia. There are a few interleaved colors, but without knowing more I’m not convinced that it can’t be explained in another way (like a mite on a host or in the gut from a carnivorous insect or even a double infection); I also can’t tell which ones are identical and which ones are just similar. But even if I should infer it from the figures, it should be reported in the results and I didn’t find it there. Maybe you are trying to state it in the subsequent sentence if one assumes that all blood feeders are the same taxa and all phloem-feeders are the same taxa, but that isn’t clear. (And at least blood feeding is a trait found in multiple diverse taxa).”

The inferences related to similar strains in distantly-related hosts is best observed in the multigene tree in figure 4 rather than the single gene trees of figures 2 and 3. For example, odonate strains are clearly interleaved between strains from other host orders. More specifically, the two Coenagrion strains have 100% identity to the Culicoides stigma strain in contrast to two other odoante (Polythore) strains where multiple SNPs are observed at all loci (See ftp file ‘BOLD\_multigene\_Rickettsia\_alignment.fas’). We thank the reviewer as this was not mentioned in the results but we have now included this on lines 209-211. Furthermore, regardless of exact MLST profiles for strains, taxa from most orders are represented in both Limoniae and Leech Torix subclades indicating a lack of grouping based on insect host taxa. The authors believe this concept is better represented in a phylogeny rather than a list of MLST profiles.

“And once again, I’m left wondering if there is a sampling bias. Are mosquitoes over-represented in the database? It some of the tables they seem over-sampled. Blood feeders and phloem feeders are often well sampled, given their important to human health and agriculture, respectively. But maybe more problematically, these results are being described but they are not clearly described in the results section. If I search for blood, I do not find any results that support this statement. When I search for phloem, there is a mention of them being found in phloem-feeding insects, but not that they are

	<p>diverse, and I have no way of assessing that as a reviewer or reader. Yet, this result for phloem-feeding is also in the abstract as a taxa that is a hot-spot. I don't see it in the figures in a way that I understand (e.g. phloem-feeding isn't annotated). Additionally, there is no assessment here that convinces me it is a hot spot; there are no statistics to suggest it is overabundant, which would be required to be a hotspot (and any such statistics would need correction for multiple testing or some sort of FDR calculation).”</p> <p>Mosquitoes are likely to be over-represented in the sequence read archive but as mentioned already in lines 142-144, a single dataset per species was extracted for analysis to negate oversampling of the same species. Although certain genera may still be oversampled, the only instance of mosquito Rickettsia being detected is in the Anopheles plumbeus population of the targeted screen. With regards to phloem-feeding insect strains, psyllids and other phloem-feeders are present in both Limoniae and Leech subclades suggesting again that strains are diverse within similar lifestyles. This is best seen in Figure 6 where both phloem-feeding and blood-feeding are annotated. The common patterns of infection in phloem-feeding bugs and blood-feeders are also already mentioned in lines 296-302 of the results. We agree that the common patterns or 'hot-spots' found in our data should come with caveats and we have now included this in the limitations part of the discussion where we clarify that common patterns of infection refer specifically to our datasets which although extensive, have some biases and may not completely represent Torix Rickettsia infection in nature (lines 358-371).</p> <p>““as previously described” Is this in a different manuscript, or earlier in the manuscript? I suspect this means earlier in the paper, but it needs to be clear. Was the same alignment method used with MAFFT and Gblocks? Same for ModelFinder? If it is and all ML trees were inferred the same way, I would recommend that you have a methods section that describes this one for all alignments, maybe concluding with what is different (like the model).”</p> <p>This refers to methods described earlier in the manuscript and yes, the same methods were used for all ML trees. This suggestion is welcomed and has been included in lines 492-493.</p> <p>“I've outlined some examples where the statements don't reflect what is presented in the results and the limitations of a secondary data analysis. But they are actually more numerous and pervasive than this. For instance, line 39-41 and 41-43 in the abstract have these issues. Likewise, Line 161-162 should read "Torix Rickettsia is the most common bacterial contaminant sequence currently in BOLD, a major barcoding project". This change reflects that this only holds for what has been barcoded thus far, and the issues with the fact you need both failed host amplification and successful bacterial amplification, and that the biodiversity represented in such projects have their own biases. It is so pervasive, I am not sure I found all the instances. Honestly, I think the paper would really benefit from a large clearly labeled section that more explicitly deals with all the limitations of the study, so others do not misconstrue the results for years into the future. It would make the paper much stronger and definitely more rigorous.”</p> <p>With regard to line 39-41 describing how our targeted PCR data supports the aquatic hotspot hypothesis we refer the reviewer to our response to the 'multiple testing' above. For lines 41-43, we have changed this sentence to include the caveat that this applies only to arthropod genome projects: “Furthermore, the analysis of 1,341 Sequence Read Archive (SRA) deposits indicates Torix infections represent a significant proportion of all Rickettsia symbioses found in arthropod genome projects.” We have also changed lines 161-162 to reflect a similar caveat for the BOLD data as suggested by the reviewer. As previously mentioned, we have now included a specific section in the discussion detailing the limitations of our datasets (lines 358-388).</p>
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No

<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	<p>Yes</p>
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>

1           **Torix *Rickettsia* are widespread in arthropods and reflect a neglected**  
2   **symbiosis**

3 Jack Pilgrim<sup>1\*</sup>, Panupong Thongprem<sup>1</sup>, Helen R. Davison<sup>1</sup>, Stefanos Siozios<sup>1</sup>, Matthew Baylis<sup>1,2</sup>,  
4 Evgeny V. Zakharov<sup>3</sup>, Sujeevan Ratnasingham<sup>3</sup>, Jeremy R. deWaard<sup>3</sup>, Craig R. Macadam<sup>4</sup>, M.  
5 Alex Smith<sup>5</sup>, Gregory D. D. Hurst<sup>1</sup>

6  
7 1. Institute of Infection, Veterinary and Ecological Sciences, Faculty of Health and Life Sciences,  
8 University of Liverpool, Liverpool, U.K.

9 2. Health Protection Research Unit in Emerging and Zoonotic Infections, Liverpool, U.K.

10 3. Centre for Biodiversity Genomics, University of Guelph, Guelph, Ontario, Canada.

11 4. Buglife – The Invertebrate Conservation Trust, Balallan House, 24 Allan Park, Stirling, U.K.

12 5. Department of Integrative Biology, University of Guelph, Guelph, Ontario, Canada.

13

14 \* Corresponding author

15 **Emails:**

16 JP: jack.pilgrim@liverpool.ac.uk

17 PT: pongthq@liverpool.ac.uk

18 HRD: hlhdavi5@liverpool.ac.uk

19 SS: sioziOSS@liverpool.ac.uk

20 MB: baylism@liverpool.ac.uk

21 EVZ: zakharov@uoguelph.ca

22 SR: sratnasi@uoguelph.ca

23 JRD: dewaardj@uoguelph.ca



24 CRM: [craig.macadam@buglife.org.uk](mailto:craig.macadam@buglife.org.uk)

25 MAS: [salex@uoguelph.ca](mailto:salex@uoguelph.ca)

26 GDDH: [ghurst@liverpool.ac.uk](mailto:ghurst@liverpool.ac.uk)

27

28 ORCID IDs:

29 Jack Pilgrim: 0000-0002-2941-1482; Panupong Thongprem: 0000-0001-6542-235X; Helen

30 Davison: 0000-0002-4302-5756; Stefanos Siozios: 0000-0002-1104-7061; Matthew Baylis;

31 Sujeevan Ratnasingham: 0000-0002-3694-0171; Gregory Hurst: 0000-0002-7163-7784

32



33 **Abstract**

34 **Background:** *Rickettsia* are intracellular bacteria best known as the causative agents of human  
35 and animal diseases. Although these medically important *Rickettsia* are often transmitted via  
36 haematophagous arthropods, other *Rickettsia*, such as those in the Torix group, appear to  
37 reside exclusively in invertebrates and protists with no secondary vertebrate host.  
38 Importantly, little is known about the diversity or host range of Torix group *Rickettsia*.

39 **Results:** This study describes the serendipitous discovery of *Rickettsia* amplicons in the  
40 Barcode of Life Data System (BOLD), a sequence database specifically designed for the  
41 curation of mtDNA barcodes. Out of 184,585 barcode sequences analysed, *Rickettsia* is  
42 observed in approximately 0.41% of barcode submissions and is more likely to be found than  
43 *Wolbachia* (0.17%). The Torix group of *Rickettsia* are shown to account for 95% of all  
44 unintended amplifications from the genus. A further targeted PCR screen of 1,612 individuals  
45 from 169 terrestrial and aquatic invertebrate species identified mostly Torix strains and  
46 supports the 'aquatic hot spot' hypothesis for Torix infection. Furthermore, the analysis of  
47 1,341 Sequence Read Archive (SRA) deposits indicates Torix infections represent a significant  
48 proportion of all *Rickettsia* symbioses found in arthropod genome projects.

49 **Conclusions:** This study supports a previous hypothesis which suggests Torix *Rickettsia* are  
50 overrepresented in aquatic insects. In addition, multiple methods reveal further putative hot  
51 spots of Torix *Rickettsia* infection; including in phloem-feeding bugs, parasitoid wasps, spiders,  
52 and vectors of disease. The unknown host effects and transmission strategies of these  
53 endosymbionts make these newly discovered associations important to inform future  
54 directions of investigation involving the understudied Torix *Rickettsia*.

55 **Keywords:** *Rickettsia*; symbiosis; arthropods; endosymbiont; DNA barcoding

## 56 **Background**

57 It is now widely recognized that animals live in a microbial world, and that many aspects of  
58 animal biology, ecology and evolution are a product of their symbioses with microorganisms  
59 [1]. In invertebrates, these symbioses may be particularly intimate, and involve transmission  
60 of the microbe from parent to offspring [2]. The alignment of host reproduction with symbiont  
61 transmission produces a correlation between the fitness interests of the parties, reflected in  
62 symbionts evolving to play a number of physiological roles within the host, from defence [3,4]  
63 through to core anabolic and digestive functions [5,6]. However, the maternal inheritance of  
64 these microbes has led to the retention of parasitic phenotypes associated with distortion of  
65 reproduction, with symbiont phenotypes including biases towards daughter production and  
66 cytoplasmic incompatibility [7]. These diverse individual impacts alter the ecology and  
67 evolution of the host, in terms of diet, dynamics of interaction with natural enemies, sexual  
68 selection and speciation.

69

70 Heritable symbioses have evolved on multiple occasions amongst microbial taxa. In some  
71 cases, the microbial lineage is limited to a single clade of related animal hosts, such as  
72 *Buchnera* in aphids [8]. In other cases, particular heritable microbes are found across a wide  
73 range of arthropod species. *Wolbachia* represents the most common associate, considered to  
74 infect nearly half of all species [9], and this commonness is a function in part of the ability of  
75 *Wolbachia* to transfer to a broad range of new host species and spread within them (host shift  
76 events) [10]. Aside *Wolbachia*, other microbes are found commonly as heritable symbionts of  
77 arthropod hosts [11]. *Cardinium* and *Rickettsia*, for instance, have been estimated at being  
78 present in 13-55% and 20-42% of terrestrial arthropod species respectively [12].

79

80 In this paper, we address the diversity and commonness of symbioses between *Rickettsia* and  
81 arthropods. The *Rickettsia* have increasingly been recognized as a genus of bacteria with  
82 diverse interactions with arthropods [13,14]. First discovered as the agents underlying several  
83 diseases of humans vectored by haematophagous arthropods [15,16], our understanding of  
84 the group changed in the 1990s with the recognition that *Rickettsia* were commonly  
85 arthropod symbionts [17,18]. *Rickettsia* were recognized first as male-killing reproductive  
86 parasites [17,19] and then later as beneficial partners [3,20,21].

87

88 Following this extension of our understanding of *Rickettsia*-arthropod interactions, a new  
89 clade of *Rickettsia* was discovered from work in *Torix* leeches [22,23]. This clade was sister to  
90 all other *Rickettsia* genera and contained two subgroups (Leech and Limoniae [24]), with no  
91 evidence to date of any strain having a vertebrate pathogen phase. The host range for *Torix*  
92 *Rickettsia* is broader than that for other members of the genus, going beyond arthropods to  
93 include amoeba hosts [25,26]. Targeted PCR based screening have revealed *Torix* group  
94 *Rickettsia* as particularly common in three groups with aquatic association: *Culicoides* biting  
95 midges, deronectid beetles and odonates [24,27,28]. However, some previous hypothesis-  
96 free PCR screens that aimed to detect *Rickettsia* in arthropods have likely missed these  
97 symbioses, due to divergence of the marker sequence and mismatch with the primers [29].

98

99 During our previous work on *Torix Rickettsia* in biting midges [27], we became aware of the  
100 presence of *Rickettsia* cytochrome *c* oxidase I (*COI*) sequences deposited in GenBank that  
101 derived from studies where the intended target of amplification/sequencing was

102 mitochondrial *COI*. These deposits derived from studies using mtDNA barcoding for  
103 phylogeographic inference [30], or in barcoding based species identification approaches  
104 [31,32]. Non-target amplification of *Rickettsia COI* using mitochondrial *COI* barcoding primers  
105 has been reported in spiders [31,32] and freshwater amphipods [30,33]. Furthermore, we  
106 have noted two cases in our lab where amplicons obtained for mtDNA barcoding of an  
107 arthropod have, on sequence analysis, revealed *Rickettsia COI* amplification (Belli group  
108 *Rickettsia* from Collembola, and Torix group *Rickettsia* from *Cimex lectularius* bedbugs).  
109 Previous work had established barcoding approaches may amplify *COI* from *Wolbachia*  
110 symbionts [34], and the data above indicated that non-target *Rickettsia COI* may be likewise  
111 amplified during this PCR amplification for mitochondrial *COI*.

112

113 In this paper, we use three approaches to reveal the diversity and commonness of Torix  
114 *Rickettsia* in arthropods. First, we probed a bin from the Barcode of Life Data System (BOLD  
115 [35]), containing non-target *COI* sequences, for *Rickettsia* amplicons and then used the DNA  
116 extracts from these projects to define the diversity of *Rickettsia* observed using a multilocus  
117 approach. Second, we screened DNA extracts from multiple individuals from 169 invertebrate  
118 species for *Rickettsia* presence to determine the distribution of the symbiont in both  
119 terrestrial and aquatic biomes. Finally, we used bioinformatic approaches to examine the  
120 Sequence Read Archive (SRA) depositions for one individual from 1,341 arthropod species for  
121 the presence of *Rickettsia* and used this as a means of estimating the relative balance of Torix  
122 group to other *Rickettsia* within symbioses.

123

124

## 125 **Data Description**

### 126 *Barcode of Life Data System (BOLD)*

127 While searching the Barcode of Life Data System (BOLD), a depository of >8 million *COI* mtDNA  
128 sequences, hundreds of hits were observed with high sequence similarity to Torix group  
129 *Rickettsia*. To investigate the diversity and host distribution of these non-target amplicons,  
130 access was permitted to analyse *COI* barcoding data deriving from a BOLD screening project  
131 totaling 184,585 arthropod specimens (including individuals where barcoding had failed) from  
132 21 countries and collected between 2010 and 2014. *COI* sequences provided by BOLD were  
133 generally derived from DNA extracts created from somatic tissues (legs are often used in order  
134 to retain most of the specimen for further analyses if necessary), but also rarely included  
135 abdominal tissues. The first dataset made available [36] included 3,817 specimens containing  
136 sequences not matching initial morphological assignment (and likely to contain contaminant  
137 sequences). The second dataset included 55,366 specimens judged to not contain non-target  
138 amplicons [37]. A remaining 125,402 specimens were not made available, and the 55,366  
139 subsample was used as a representative sample from which the contaminants had originated  
140 (Figure 1). The protocols for data collection, data curation and quality control of submitted  
141 BOLD samples is described by Ratnasingham & Hebert [38].

142

### 143 *Sequence Read Archive (SRA)*

144 Further insights into the balance of *Rickettsia* groups within arthropod symbioses were  
145 obtained through searching for *Rickettsia* presence in Illumina datasets associated with  
146 arthropod whole genome sequence (WGS) projects in the SRA (60,409 records as of the 20th  
147 May 2019). To reduce the bias from over-represented laboratory model species (e.g.

148 *Drosophila* spp., *Anopheles* spp.) a single dataset per species was examined, and where  
149 multiple data sets existed for a species, that with the largest read count was retained. The  
150 resultant dataset [39], representing 1,341 arthropod species, was then screened with  
151 phyloFlash [40] which finds, extracts and identifies SSU rRNA sequences.

152

### 153 *Targeted screen of aquatic and terrestrial arthropods*

154 Both the BOLD and SRA datasets have inherent biases which make them unsuitable to assess  
155 whether *Torix Rickettsia* are more common in aquatic or terrestrial biomes. For example, most  
156 SRA submissions are from lab-reared terrestrial insects. Likewise, a majority of the BOLD  
157 specimens containing *Rickettsia* have limited taxonomic and ecological information, by virtue  
158 of not returning an mtDNA *COI* sequence. Therefore, a targeted PCR screen of 1,612  
159 individuals from 169 species was undertaken (Table 1) using primers which hybridise with all  
160 known clades of *Rickettsia* [27]. Within this, we included a range of both aquatic and terrestrial  
161 taxa, to investigate if the previous work highlighting particular aquatic taxa as hot spots for  
162 *Rickettsia* symbiosis (water beetles, biting midges, damselflies) reflects a wider higher  
163 incidence in species from this habitat.

164

## 165 **Analyses**

166 *Torix Rickettsia* is the most common bacterial contaminant sequence currently in BOLD, a  
167 major barcoding project

168 Out of 3,817 sequences considered as not matching initial morphological assignment, 1,126  
169 of these were deemed by BOLD to be bacterial in origin (Figure 1, [36]). The taxonomic  
170 classification tool, Kaiju, further supported bacterial designation for all sequences except one

171 (Additional file 1), although this was later confirmed as *Rickettsia* through phylogenetic  
172 placement. Phylogenetic placement further confirmed the correct designation of bacterial  
173 sequences (Figure 2 and Additional file 2). The dominant genus was *Rickettsia* with 753  
174 (66.9%) amplifications, compared to *Wolbachia* with 306 (27.2%). Of the remaining 67 non-  
175 target sequences, 14 formed a monophyletic group with other Anaplasmataceae and 48  
176 clustered with the order Legionellales, with 5 sequences remaining undesignated. When  
177 considering the 184,585 specimens in the total project, this analysis gave an overall *Rickettsia*  
178 and *Wolbachia* frequency of 0.41% and 0.17% respectively within the dataset. Through later  
179 access to the 55,366 representative data subset from where the contaminants originated, a  
180 further 245 unique bacteria contaminants were also detected by Kaiju (possibly missed by  
181 BOLD's automated contaminant filtering system) (Additional file 1). This additional finding  
182 suggests these frequencies are conservative estimates.

183

184 BOLD *Rickettsia* contaminants were dominated by amplicons from the Torix group of  
185 *Rickettsia* (716/753; 95.1%) (Figure 3 and Additional file 2). The remaining 37 *Rickettsia*  
186 clustered with Transitional/Spotted Fever (n=15), Belli (n=9), Rhyzobius (n=1) groups, while 12  
187 sequences formed two unique clades. Across arthropod hosts: 292 (38.8%) were derived from  
188 Hymenoptera; 189 (25.1%) from Diptera; 177 from Hemiptera (23.5%); 41 from Psocoptera  
189 (5.4%); 40 from Coleoptera (5.3%); 7 from Arachnida (0.9%); 4 from Trichoptera (0.5%); and  
190 single cases of Thysanoptera, Diplopoda and Dermaptera (0.1% each).

191

192 We observed that two sets of *COI* primers were responsible for 99% of *Rickettsia*  
193 amplifications (Additional file 3) with a majority (89%) amplifying with the primer combination



194 C\_LepFolF/C\_LepFolR [41]. Torix *Rickettsia* *COI* showed a stronger match to these primers at  
195 the 3' end (the site responsible for efficient primer annealing) compared to *Wolbachia* and  
196 other *Rickettsia* groups. Whilst all contained a SNP at the 3' priming end of C\_LepFolR, Torix  
197 *Rickettsia* (*Rickettsia* endosymbiont of *Culicoides newsteadii*; MWZE00000000) was the only  
198 sequence to not contain a SNP at the 3' priming site of C\_LepFolF (Additional file 4).

199

#### 200 *Rickettsia* multilocus phylogenetic analysis

201 To better resolve the phylogenetic relationships between BOLD *Rickettsia* contaminants, a  
202 multilocus approach was employed on a subsample of 186 *Rickettsia*-containing samples  
203 chosen based on assorted geographic location, host order and phylogenetic placement. To  
204 this end, 2 further housekeeping genes (*16S rRNA*, *gltA*) and the antigenic *17KDa* protein gene  
205 were amplified and sequenced from the respective DNA extracts.

206

207 Overall, 135 extracts successfully amplified and gave a high-quality sequence for at least one  
208 gene. No intragenic or intergenic recombination was detected for any of the gene profiles. A  
209 phylogram, including 99 multilocus profiles containing at least 3 of the 4 *Rickettsia* genes of  
210 interest (including *COI*), allocated strains to both Limoniae and Leech subclades of the Torix  
211 group (Figure 4) and these subclades were derived from similar hosts. For example, specific  
212 families (Hemiptera: Psyllidae and Hymenoptera: Diapriidae) were present in both Leech and  
213 Limoniae groups. Furthermore, similar strains were observed between genetically dissimilar  
214 host species. For example, the *Coenagrion mercuriale* (Odonata) strain was 100% identical to  
215 the *Culicoides stigma* (Diptera) strain across all four loci. This suggests horizontal transfer of

216 the symbiont is likely to be occurring. A full list of multilocus profiles and *Rickettsia* group  
217 designation can be found in Additional file 5.

218

219 The multilocus study also provided evidence of co-infection with *Rickettsia*. During Sanger  
220 chromatogram analysis, double peaks were occasionally found at third codon sites from  
221 protein coding genes. This pattern was observed in 6/10 *Philotarsus californicus* individuals  
222 and in one member of each of the Psilidae, Sciaridae, Chironomidae and Diapriidae (Additional  
223 file 5). Where double peaks were observed, this was found consistently across markers within  
224 an individual specimen. This pattern corroborates a recent finding of double infections in  
225 Odoantes [28], suggesting co-infecting *Rickettsia* strains in hosts is a widespread phenomenon  
226 of the Torix group.

227

#### 228 *Barcoding success of Rickettsia host taxa*

229 An available subset of specimens associated with the contaminants contained 55,366 out of  
230 184,585 arthropods originally used in the overall study [37]. The three classes of Insecta  
231 (n=49,688), Arachnida (n=3,626) and Collembola (n=1,957), accounted for >99.8% of total  
232 specimens (Figure 1). Successful amplification and sequencing of *COI* was achieved in 43,246  
233 specimens (78.1%) of the DNA extracts, but when assessed at the order level success rates  
234 varied (Additional file 6). The likely explanation for this variation is taxa-specific divergence of  
235 sequences at priming sites.

236

237 The number of each taxonomic order giving at least one *Rickettsia* amplification was then  
238 calculated and adjusted based on the total number of specimens in the project to allow for a

239 frequency estimate. Overall, Hymenoptera, Diptera and Hemiptera were the three taxa most  
240 likely to be associated with *Rickettsia COI* amplification (87.4%). Similarly, on assessment of a  
241 subsample from the project where the contaminants originated, a majority (77.7%) of the  
242 dataset were also accounted for by these three orders. After adjusting the frequency to take  
243 into account the number of inaccessible specimens, Trichoptera (2.45%), Dermaptera (1.89%)  
244 and Psocodea (1.67%) were the most likely taxa to give an inadvertent *Rickettsia* amplification.  
245 Whilst Hemiptera and Diptera had a similar estimated frequency of *Rickettsia* amplification  
246 (0.58% and 0.56%), Hemiptera were much more likely to fail to barcode (67.2% vs 93.3%),  
247 suggesting dipteran *Rickettsia* infection in BOLD specimens is likely to be higher than that of  
248 hemipterans, as a barcoding failure is necessary to amplify non-target bacteria *COI*. Attempts  
249 to re-barcode 186 *Rickettsia*-containing DNA extracts of interest from BOLD resulted in 90  
250 successful arthropod host barcodes (Additional file 5).

251

#### 252 *Targeted Rickettsia PCR screen and statistical comparison of terrestrial vs aquatic insects*

253 From the targeted screen of 169 invertebrate species, a total of 19 *Rickettsia* were discovered  
254 from both aquatic and terrestrial pools, with 17 of these identified as belonging to the Torix  
255 group. The screening of aquatic invertebrates revealed 9 out of 57 species (15.79%) were  
256 positive in PCR assays (Table 1.1). DNA sequences confirmed that all were *Rickettsia* which lay  
257 within the Torix group (Figure 5), with the positive species deriving from 8 insect species and  
258 one mollusc. For the terrestrial invertebrates, PCR assays evidenced *Rickettsia* infection in 10  
259 out of 112 species (8.93%) with a mix of insect and spider hosts (4 and 6 species respectively,  
260 Table 1.2). *Rickettsia* from 8 host species (2 insects and 6 spiders) were identified as Torix

261 *Rickettsia* (8 of 112 species, 7.14%), while the other two host species carried *Rickettsia* from  
262 the Rhyzobius and Belli groups (Figure 5).

263

264 To reduce taxonomic hot spot biases (particularly from spiders), we compared the incidence  
265 of *Rickettsia* infection in aquatic vs terrestrial insects. Fisher's exact test analysis rejected the  
266 null hypothesis of equal representation, with aquatic taxa having a higher representation of  
267 species with Torix *Rickettsia* than terrestrial ( $p$ -value = 0.013, Additional file 7). Examining the  
268 phylogenetically controlled set, with three matched insect orders (Coleoptera, Diptera,  
269 Hemiptera), again rejected the null hypothesis of equal representation, with aquatic taxa  
270 having a higher representation of species with Torix *Rickettsia* than terrestrial ( $p$ -value =  
271 0.025, Additional file 7). When comparing all invertebrate species from the targeted screen,  
272 no significant difference was observed in Torix *Rickettsia* incidence between terrestrial and  
273 aquatic biomes ( $p$ -value = 0.11, Additional file 7) suggesting this pattern of infection may be  
274 specific to insects.

275

276 [Insert Table 1 here]

277

278 *SRA and GenBank Rickettsia searches*

279 During the SRA search, phyloFlash flagged 29 *Rickettsia* sequences in the groups: Belli (n=10),  
280 Torix (n=8), Transitional (n=6), Rhyzobius (n=2), and Spotted Fever (n=1), with the remaining  
281 two failing to form a monophyletic clade with any group (Figure 5). In addition, Kraken  
282 identified eight *Rickettsia*-containing arthropod SRA datasets missed by phyloFlash. Two of  
283 these were from the Torix group, in phantom midge hosts (Diptera: Chaoboridae: *Mochlonyx*

284 *cinctipes* and *Chaoborus trivitattus*), with the remaining six placed in Belli and Spotted Fever  
285 groups [39].

286

287 phyloFlash was also used to retrieve 18S rRNA (eukaryotic) sequences which could potentially  
288 account for the *Rickettsia* observed in SRA datasets (e.g. through parasitisms or ingestion of  
289 *Rickettsia*-infected protists). Out of the 29 datasets analysed by phyloFlash, only one  
290 (SRR6313831) revealed an assembled 18S rRNA sequence aligned to a parasitoid wasp  
291 (*Hadrotrichodes waukheon*). Although reads aligned to protists were also present in 19/29  
292 datasets flagged by phyloFlash, the read depth for protists was much lower than the number  
293 of *Rickettsia* reads [39]. This suggests that *Rickettsia*-infected protists are unlikely to account  
294 for the positives observed in the SRA datasets.

295

296 The search of GenBank revealed 11 deposits ascribed to host mtDNA that were in fact Torix  
297 *Rickettsia* sequences (Additional files 8 and 9).

298

### 299 *The hidden host diversity of Torix Rickettsia*

300 Overall, putative novel Torix hosts detected from all screening methods included taxa from  
301 the orders Dermoptera, Gastropoda, Trichoptera and Trombidiformes. Additionally, new  
302 Torix-associated families, genera and species were identified. These included  
303 haematophagous flies (*Simulium aureum*; *Anopheles plumbeus*; *Protocalliphora azurea*;  
304 Tabanidae), several parasitoid wasp families (e.g. Ceraphronidae; Diapriidae; Mymaridae),  
305 forest detritivores (e.g. Sciaridae; Mycetophilidae; Staphylinidae) and phloem-feeding bugs  
306 (Psyllidae; Ricaniidae). Feeding habits such as phloem-feeding, predation, detritivory or

307 haematophagy were not correlated with any particular Torix *Rickettsia* subclade (Figure 6).  
308 Furthermore, parasitoid and aquatic lifestyles were seen across the phylogeny. All newly  
309 discovered putative Torix *Rickettsia* host taxa are described in Table 2, alongside previously  
310 discovered hosts in order to give an up to date overview of Torix-associated taxa.

311

312

313 [Insert Table 2 here]

314

315

## 316 **Discussion**

317 Symbiotic interactions between hosts and microbes are important drivers of host phenotype,  
318 with symbionts both contributing to, and degrading, host performance. Heritable microbes  
319 are particularly important contributors to arthropod biology, with marked attention focused  
320 on *Wolbachia*, the most common associate [9]. Members of the Rickettsiales, like *Wolbachia*,  
321 share an evolutionary history with mitochondria [42], such that a previous screen of BOLD  
322 submissions of mtDNA submissions observed *Wolbachia* as the main bacterial contaminant  
323 associated with DNA barcoding [34]. However, our screen found that *Rickettsia* amplicons  
324 were more commonly found in BOLD deposits compared to *Wolbachia* (0.41% vs 0.17% of  
325 deposits). Furthermore, Torix group *Rickettsia* were overrepresented in barcode  
326 misamplifications (95%) when compared to other groups within the genus. A comparison of  
327 the most commonly used barcoding primers to *Wolbachia* and *Rickettsia* genomes suggest  
328 homology of the forward primer 3' end was likely responsible for this bias towards Torix

329 *Rickettsia* amplification. To gain a clearer understanding of the relative balance of Torix group  
330 to other *Rickettsia* within symbioses and habitats, a targeted screen and bioinformatic  
331 approach was also undertaken. Through these three screens, a broad range of host diversity  
332 associated with Torix *Rickettsia* was uncovered.

333

334 As the *in silico* and empirical evidence suggests *Rickettsia COI* amplification is not uncommon  
335 [31–33], why has this phenomenon not been described more widely before? The previous  
336 large-scale non-target *COI* study using BOLD submissions [34], revealed only *Wolbachia* hits.  
337 This screen involved comparison to a *Wolbachia*-specific reference library and was thus likely  
338 to miss *Rickettsia*. Additionally, there has been a lack of Torix *Rickettsia COI* homologues to  
339 compare barcodes to until recently, where a multilocus identification system, including *COI*  
340 was devised [27]. Indeed, out of the non-target *COI* dataset received in this study, some of the  
341 *Rickettsia* contaminants were tentatively described by BOLD as *Wolbachia* due to the previous  
342 absence of publicly available *Rickettsia COI* to compare.

343

344 Although *Rickettsia* will only interfere with barcoding in a minority of cases (~0.4%), it is likely  
345 that alternate screening primers for some studies will need to be considered. In a  
346 demonstration of how unintended *Rickettsia* amplifications can affect phylogeographic  
347 studies relying on DNA barcoding, a *Rickettsia COI* was conflated with the mtDNA *COI* of a  
348 species of freshwater amphipod, *Paracalliope fluvialis* [30]. Subsequently, supposed unique  
349 mtDNA haplotypes were allocated to a particular collection site, whereas this merely  
350 demonstrated the presence of Torix *Rickettsia* in host individuals in this lake. Contrastingly,



351 non-target *Rickettsia* amplification can also allow for the elucidation of a novel host range of  
352 the symbiont [31–33] and this has been exemplified with our probing of BOLD.

353

354 Previously, several host orders have been associated with Torix *Rickettsia*, including Araneae,  
355 Coleoptera, Diptera, Hemiptera and Odonata [24,28,43–45]. Newly uncovered putative host  
356 orders from this study include Dermaptera, Gastropoda, Trichoptera and Trombidiformes  
357 (Table 2). These data emphasise the broad host range of Torix *Rickettsia* across arthropods  
358 and invertebrates, with two additional cases from nucleariid amoebae [25,26]. This host range  
359 is complementary to *Rickettsia*'s sister genus '*Candidatus Megaira*' (formally the Hydra group  
360 of *Rickettsia*) which are present in multiple unicellular eukaryote families, and in a few  
361 invertebrates like *Hydra* [46].

362

363 Despite the extensive sampling and multiple screening strategies employed in this project,  
364 caution must be taken when interpreting to what extent the Torix *Rickettsia* hosts identified  
365 are representative of *Rickettsia* hosts in nature. Both BOLD and SRA components of the project  
366 rely on secondary data which come with sampling and methodological biases. For example,  
367 most SRA submissions are from lab-reared terrestrial insects and it can be argued that the  
368 high number of Belli *Rickettsia* infections discovered from arthropod genome projects  
369 (compared to the targeted screen which contains multiple aquatic insect species) could be  
370 due to this sampling bias. Likewise, the over-representation of Torix *Rickettsia* from BOLD is  
371 likely due to an amplification bias as a result of higher primer site homology to that particular  
372 group from commonly used barcoding primer sets. Subsequently, the common patterns of  
373 infection (or 'hot spots') found in this study are identified as such with these provisos in mind.

374 To counteract these biases and to give a more nuanced and holistic view of *Torix Rickettsia*  
375 ecology, a targeted screen was also included to ensure this study was not over-reliant on  
376 secondary data.

377

378 Further caution needs to be taken when interpreting what these newly found associations  
379 mean, as mere presence of *Rickettsia* DNA does not definitively indicate an endosymbiotic  
380 association. For example, bacterial DNA integrations into the host nuclear genome have been  
381 widely reported [47]. Although none of the protein-coding genes sequenced in this study  
382 showed signs of a frameshift, suggesting a lack of pseudogenization that is often typical of a  
383 nuclear insertion, this still does not rule out this phenomenon entirely. Furthermore,  
384 parasitism or ingestion of symbiont-infected biota (e.g. protists) could also result in bacteria  
385 detection [48–50]. Whilst protist reads were found in some datasets, these were usually at a  
386 much lower depth compared to the symbiont [39]. In one of the few instances where protist  
387 reads were greater than *Rickettsia* (Dataset SRR5298327), this was from our own previous  
388 study where a true endosymbiosis between insect and symbiont was confirmed through FISH  
389 imaging [27]. Similarly, although an 18S sequence aligned to a parasitoid wasp was observed  
390 in the SRA dataset from *Bemisia tabaci* (SRR6313831), previous work has also demonstrated  
391 a true endosymbiosis between *B. tabaci* and *Torix Rickettsia* [51]. Overall, these data suggest  
392 that detecting contamination from *Rickettsia*-infected taxa such as protists and parasitoid  
393 wasps is uncommon within our study.

394

395 Model-based estimation techniques suggest *Rickettsia* are present in between 20-42% of  
396 terrestrial arthropod species [12]. However, the targeted PCR screen in this study gave an

397 estimated species prevalence of 8.9% for terrestrial species. This discrepancy is likely due to  
398 targeted screens often underestimating the incidence of symbiont hosts due to various  
399 methodological biases including small within-species sample sizes (missing low-prevalence  
400 infections) [29]. Importantly, the inclusion and exclusion of specific ecological niches can also  
401 lead to a skewed view of *Rickettsia* symbioses. A previous review of *Rickettsia* bacterial and  
402 host diversity by Weinert et al. [13] suggested a possible (true) bias towards aquatic taxa in  
403 the Torix group. In accordance with this, our targeted screen demonstrated Torix *Rickettsia*  
404 infections were more prevalent in aquatic insect species compared to terrestrial (although this  
405 is likely not the case for invertebrates in general due to a Torix *Rickettsia* hot spot in spiders).  
406 The observed over-representation of Torix group *Rickettsia* (17/19 strains) in our targeted  
407 screen contrasts with Weinert's findings which show a predominance of Belli infections and is  
408 likely due to the latter study's near absence of aquatic insects and spiders within the samples  
409 screened. Our additional use of a bioinformatics approach based on the SRA appears to  
410 corroborate targeted screen data indicating that Belli and Torix are two of the most common  
411 *Rickettsia* groups among arthropods. Overall, these multiple screening methods suggest Torix  
412 *Rickettsia* are more widespread than previously thought and their biological significance  
413 underestimated.

414

415 Previous studies have used either one or two markers to identify the relatedness of strains  
416 found in distinct hosts. In this study, we use the multilocus approach developed in Pilgrim et  
417 al. [27] to understand the affiliation of Torix *Rickettsia* from diverse invertebrate hosts. Our  
418 analysis of Torix strains indicates that closely related strains are found in distantly related taxa.  
419 Closely related *Rickettsia* are also found in putative hosts from different niches and habitats –

420 for instance, the *Rickettsia* strains found in terrestrial blood feeders do not lie in a single clade,  
421 but rather are allied to strains found in non-blood feeding host species. Likewise, strains in  
422 phloem-feeding insects are diverse rather than commonly shared.

423

424 The distribution of Torix *Rickettsia* across a broad host range suggests host shifts are occurring  
425 between distantly related taxa. It is notable that parasitoid wasps are commonly infected with  
426 *Rickettsia* and have been associated with enabling symbiont host shifts [48]. Aside from  
427 endoparasitoids, it is also possible that plant-feeding can allow for endosymbiont horizontal  
428 transmission [52,53]. For example, *Rickettsia* horizontal transmission has been demonstrated  
429 in *Bemisia* whiteflies infected by phloem-feeding [52,54]. Finally, ectoparasites like the Torix-  
430 infected water mites of the Calyptostomatidae family, could also play a role in establishing  
431 novel *Rickettsia*-host associations, as feeding by mites has been observed to lead to host shifts  
432 for other endosymbiont taxa [55]. Indeed, if multiple horizontal transmission paths do exist,  
433 this could account for the diverse plethora of infected taxa, as well as arthropods identified in  
434 this study which harbour more than one strain of symbiont [56].

435

436 The finding that Torix *Rickettsia* are associated with a broad range of invertebrates leads to  
437 an obvious question: what is the impact and importance of these symbiotic associations?  
438 Previous work has established Torix *Rickettsia* represent heritable symbionts and it is likely  
439 that this is true generally. There have, however, been few studies on their impact on the host.  
440 In the earliest studies [22,23], *Torix* spp. leeches infected with *Rickettsia* were observed to be  
441 substantially larger than their uninfected counterparts. Since then, the only observation of  
442 note, pertaining to the Torix group, is the reduced ballooning (dispersal) behaviour observed

443 in infected *Erigone atra* money spiders [57]. Overall, the incongruencies in host and Torix  
444 *Rickettsia* phylogenies (suggesting a lack of co-speciation and obligate mutualism), along with  
445 the lack of observed sex bias in carrying the symbiont, indicate facultative benefits are the  
446 most likely symbiotic relationship [29]. However, *Rickettsia* induction of thelytokous  
447 parthenogenesis (observed in Belli *Rickettsia* [58,59]) should not be discounted in Torix  
448 infected parasitoid wasps identified in this study. To add to the challenge of understanding  
449 Torix *Rickettsia* symbioses, the challenges of laboratory rearing of many Torix *Rickettsia* hosts  
450 has led to difficulties in identifying model systems to work with. However, the large expansion  
451 of our Torix group host knowledge can now allow for a focus on cultivatable hosts (e.g phloem-  
452 feeding bugs).

453

454 To conclude, we have shown that large-scale DNA barcoding initiatives of arthropods can  
455 include non-target amplification of Torix *Rickettsia*. By examining these non-target sequences,  
456 alongside a targeted screen and SRA search, we have uncovered numerous previously  
457 undetected putative host associations. Our findings lay bare multiple new avenues of inquiry  
458 for Torix *Rickettsia* symbioses.

459

## 460 **Potential Implications**

461 A particularly important group for future study of Torix *Rickettsia* interactions are  
462 haematophagous host species. Our discovery of *Rickettsia*-associated tabanid and simuliid  
463 flies, alongside *Anopheles plumbeus* mosquitoes, add to existing blood-feeders previously  
464 identified as Torix group hosts which include sand flies [60,61], fleas [62], ticks [63,64] bed  
465 bugs [65] and biting midges [27]. Some *Rickettsia* strains are known to be transmitted to

466 vertebrates via haematophagy [66]. However, there is no evidence to date for vertebrate  
467 pathogenic potential for the Torix group. Despite this, Torix *Rickettsia* could still play a  
468 significant role in the ecology of vectors of disease. A key avenue of research is whether these  
469 endosymbionts alter vectorial capacity, as found for other associations [67]. In contrast to the  
470 widely reported virus blocking phenotype observed in *Wolbachia*-infected vectors [68,69],  
471 Torix *Rickettsia* has recently been associated with a virus potentiating effect in *Bemisia* white  
472 flies vectoring Tomato yellow leaf curl virus [70]. Additionally, we uncovered a *Rickettsia*-  
473 infected psyllid (*Cacopsylla melanoneura*) which is a vector of *Phytoplasma mali* (apple  
474 proliferation) [71]. Thus, the question of Torix *Rickettsia* vector-competence effects is clearly  
475 of widespread relevance and deserves further attention.

476

## 477 **Methods**

### 478 **a) Interrogation of the Barcode of Life Data System (BOLD)**

#### 479 *Assessment of non-target microbe amplicons*

480 BOLD data curation involves identifying non-target *COI* sequences from common  
481 contaminants (e.g. human and bacteria) or erroneous morphological identifications [38]. The  
482 designation of bacterial contaminants by BOLD, from a dataset containing 3,817 non-target  
483 sequences [36], was confirmed by the taxonomic classification program, Kaiju, using default  
484 parameters [72]. Sequences were then placed phylogenetically to refine taxonomy further.  
485 To this end, barcodes confirmed as microbial sequences were aligned using the “L-INS-I”  
486 algorithm in MAFFT v7.4 (RRID:SCR\_011811) [73]. Gblocks (RRID:SCR\_015945) [74] was then  
487 used to exclude areas of the alignment with excessive gaps or poor alignment using ‘options  
488 for a less stringent selection’; the inclusion of some missing data in alignments was allowed as

489 missing characters does not often affect phylogenetic resolution for taxa with complete data  
490 [75]. ModelFinder [76] then determined the TIM3+F+I+G4 model to be used after selection  
491 based on default “auto” parameters using the Bayesian information criteria. A maximum  
492 likelihood (ML) phylogeny was then estimated with IQTree [77] using an alignment of 561  
493 nucleotides and 1000 ultrafast bootstraps [78]. The Rickettsiales genera *Anaplasma*,  
494 *Rickettsia*, *Orientia* and *Wolbachia* (Supergroups A, B, E and F), as well as the Legionellales  
495 genera *Legionella* and *Rickettsiella*, were included in the analysis as references (as suggested  
496 by Kaiju). Finally, both phylogram and cladogram trees (the latter for ease of presentation)  
497 were drawn and annotated based on host taxa (order) using the EvolView [79] online tree  
498 annotation and visualisation tools. Subsequent phylogenetic workflows detailed below follow  
499 this method with the exception being the chosen models by Modelfinder.

500

501 A determining factor for non-target amplification of bacteria is primer site matching to  
502 microbial associates. Subsequently, pairwise homology of the primer set predominantly used  
503 for BOLD barcode screening was compared to *Rickettsia* and *Wolbachia COI* genes.

504

#### 505 *Further phylogenetic analysis*

506 *COI* sequence alone provides an impression of the frequency with which *Rickettsia* associates  
507 are found in barcoding studies. However, they have limited value in describing the diversity of  
508 the *Rickettsia* found. To provide further insight into the diversity of *Rickettsia* using a  
509 multilocus approach, we obtained 186 DNA extracts from the archive at the Centre for  
510 Biodiversity Genomics (University of Guelph, Canada) that had provided *Rickettsia* amplicons  
511 in the previous screen. DNA extracts were chosen based on assorted geographic location, host



512 order and phylogenetic placement. Multilocus PCR screening and phylogenetic analysis of  
513 *Rickettsia* was then completed, using the methodology in Pilgrim et al. which utilised primers  
514 conserved across all known clades of the *Rickettsia* genus [27]. However, slight variations  
515 include the exclusion of the *atpA* gene due to observed recombination at this locus.  
516 Furthermore, the amplification conditions for the *17KDa* locus was changed because a Torix  
517 *Rickettsia* reference DNA extract (Host: *Simulium aureum*) failed to amplify with the primer  
518 set Ri\_17KD\_F/ Ri\_17KD\_R from Pilgrim et al. [27]. Subsequently, a *17KDa* alignment from  
519 genomes spanning the Spotted fever, Typhus, Transitional, Belli, Limoniae groups, and the  
520 genus '*Candidatus Megaira*' was generated to design a new set of primers using the online  
521 tool PriFi [80].

522

523 Once multilocus profiles of the *Rickettsia* had been established, we tested for recombination  
524 within and between loci using RDP v4 (Recombination Detection Program, RRID:SCR\_018537)  
525 [81] using the MaxChi, RDP, Chimaera, Bootscan and GENECONV algorithms with the following  
526 criteria to assess a true recombination positive: a p-value of <0.001; sequences were  
527 considered linear with 1000 permutations being performed. Samples amplifying at least 3 out  
528 of 4 genes (*16S rRNA*, *17KDa*, *COI* and *gltA*) were then concatenated and their relatedness  
529 estimated using maximum likelihood as described above. The selected models used in the  
530 concatenated partition scheme [82] were as follows: *16S rRNA*: TIM3+F+R2; *17KDa*:  
531 GTR+F+I+G4; *COI*:TVM+F+I+G4; *gltA*: TVM+F+I+G4. Accession numbers for all sequences used  
532 in phylogenetic analyses can be found in Additional file 10.

533

534 *Re-barcoding Rickettsia-containing BOLD DNA extracts*

535 Aside from phylogenetic placement of these *Rickettsia*-containing samples, attempts were  
536 made to extract an mtDNA barcode from these taxa in order to identify the hosts of infected  
537 specimens. This is because morphological taxonomic classification of specimens in BOLD is  
538 usually only down to the order level before barcoding takes place. Previous non-target  
539 amplification of *Rickettsia* through DNA barcoding of arthropod DNA extracts had occurred in  
540 the bed bug *Cimex lectularius*, with a recovery of the true barcode after using the primer set  
541 C1-J-1718/HCO1490, which amplifies a shortened 455 bp sequence within the *COI* locus.  
542 Subsequently, all samples were screened using these primers or a further set of secondary *COI*  
543 primers (LCOt\_1490/ MLepR1 and LepF1/C\_ANTMR1D) if the first failed to give an adequate  
544 host barcode. All *COI* and *Rickettsia* multilocus screening primer details, including references,  
545 are available in Additional file 11.

546

547 Cycling conditions for *COI* PCRs were as follows: initial denaturation at 95°C for 5 min, followed  
548 by 35 cycles of denaturation (94°C, 30 sec), annealing (50°C, 60 sec), extension (72°C, 90 sec),  
549 and a final extension at 72°C for 7 min. *Rickettsia* and host amplicons identified by gel  
550 electrophoresis were subsequently purified enzymatically (ExoSAP) and Sanger sequenced  
551 through both strands using a BigDye® Terminator v3.1 kit (Thermo Scientific, Waltham, USA),  
552 and capillary sequenced on a 3500 xL Genetic Analyser (Applied Biosystems, Austin, USA).  
553 Forward and reverse reads were assessed in UGENE (RRID:SCR\_005579)[83] to create a  
554 consensus sequence by eye with a cut-off phred (Q) score [84] of 20. Primer regions were  
555 trimmed from barcodes before being matched to the GenBank database by BLAST based on  
556 default parameters and an e-value threshold of <1e-85. Host taxonomy was determined by a

557 barcode-based assignment of the closest BLAST hit, under the following criteria modified from  
558 Ramage et al. [50]:

559 1) Species level designation for at least 98% sequence identity.

560 2) Genus level designation for at least 95% sequence identity.

561 3) Family level designation for at least 85% sequence identity.

562 Additionally, all sequences were required to be at least >200 bp in length.

563

#### 564 *Assessment of barcoding success*

565 One of the factors determining a successful *COI* bacterial amplification is the initial failure of  
566 an extract to amplify mtDNA. Subsequently, to determine the likelihood of this event within  
567 taxa, we used the 55,366 specimen representative data subset [37] to evaluate failure rates.

568 To this end, all orders of host which gave at least one non-target *Rickettsia COI* hit were  
569 assessed. The barcoding success rate was determined as the proportion of specimens which  
570 matched initial morphotaxa assignment and were not removed after BOLD quality control  
571 [38]. As the total *Rickettsia* count was from a larger dataset than the one made available, an  
572 adjusted infection frequency for each taxon was calculated based on the representative data  
573 subset.

574

#### 575 **b) Targeted and bioinformatic *Rickettsia* screens**

##### 576 *Targeted screen of aquatic and terrestrial arthropods*

577 Overall, 1,612 individuals from 169 species, including both terrestrial (DNA extracts derived  
578 from European material, mostly from Duron et al. [11]) and aquatic invertebrates (largely  
579 acquired from the UK between 2016-2018), were screened. mtDNA *COI* amplification was

580 conducted as a control for DNA quality. Some arthropods which could not be identified down  
581 to the species level morphologically or from barcoding were referred to as 'sp.'. To investigate  
582 symbiont infection status, rickettsial-specific primers based on *gltA* and *16S rRNA* genes were  
583 used for conventional PCR screening [27], with Sanger sequences obtained from at least one  
584 specimen per *Rickettsia* positive species to identify any misamplification false positives. Newly  
585 identified hosts of interest from BOLD and targeted screens were then placed phylogenetically  
586 (see sections above) with the models TIM3+F+R2 (16S) and K3Pu+F+G4 (*gltA*) before being  
587 mapped by lifestyle and diet.

588

589 It is known that there are taxonomic hot spots for endosymbiont infection, with for instance  
590 spiders being a hot spot for a range of microbial symbionts [43]. Therefore, analyses were  
591 performed that were matched at a taxonomic level (i.e. each taxon was represented in both  
592 the aquatic and terrestrial pools). To this end, the incidence of *Torix Rickettsia* was first  
593 compared in all insects. However, within insects, there is taxon heterogeneity between  
594 aquatic and terrestrial biomes (e.g. Ephemeroptera, Plecoptera in aquatic only, Lepidoptera  
595 in terrestrial only). The analysis was therefore narrowed to match insect orders present in  
596 both the aquatic and terrestrial community. Three insect orders, Hemiptera, Diptera and  
597 Coleoptera, fulfilled this criterion with good representation from each biome. For each case,  
598 the ratios of the infected:non-infected species between aquatic and terrestrial communities  
599 were compared in a Fisher's exact test with a *p*-value significance level of  $\leq 0.05$ .

600

601 *Search of the Sequence Read Archive (SRA) and GenBank*

602 The SRA dataset [39] containing one individual from 1,341 arthropod species was screened  
603 with phyloFlash [40] using default parameters, which finds, extracts and identifies SSU rRNA  
604 sequences. Reconstructed full *16S rRNA* sequences affiliated to *Rickettsia* were extracted and  
605 compared to sequences derived from the targeted screen phylogenetically (see sections  
606 above) to assess group representation within the genus. The microbial composition of all SRA  
607 datasets that did not result in a reconstructed *Rickettsia 16S rRNA* with phyloFlash were re-  
608 evaluated using Kraken2 [85], a k-mer based taxonomic classifier for short DNA sequences. A  
609 cut-off of at least 40k reads assigned to *Rickettsia* taxa was applied for reporting potential  
610 infections (theoretical genome coverage of  $\sim 1 - 4X$  assuming an average genome size of  
611  $\sim 1.5\text{Mb}$ ). As *Rickettsia*-infected protists and parasitoids have previously been reported  
612 [25,26,59], phyloFlash was also used to identify reads aligned to these taxa to account for  
613 potential positives attributed to ingested protists or parasitisms.

614

615 We also examined GenBank for *Rickettsia* sequences deposited as invertebrate *COI* barcodes.  
616 To this end, a BLAST search of Torix *Rickettsia COI* sequences from previous studies [27,32]  
617 was conducted on the 29<sup>th</sup> June 2020. Sequences were putatively considered belonging to the  
618 Torix group if their similarity was  $>90\%$  and subsequently confirmed phylogenetically as  
619 described above with the HKY+F+G4 model.

620

621 **Table 1.1.** Targeted *Rickettsia* screen of aquatic/semiaquatic invertebrates.

622

Aquatic/Semiaquatic invertebrate group	Species	Location	Year	No. tested	No positive
	<i>Baetis muticus</i>	Stirling, Scotland, UK	2017	3	0
	<i>Baetis rhodani</i>	Stirling, Scotland, UK	2017	3	0
	<i>Cloeon dipterum</i>	Cheshire, UK	2016	3	0
	<i>Ecdyonurus</i> sp.1	Stirling, Scotland, UK	2017	5	0

Ephemeroptera	<i>Ecdyonurus</i> sp.2	Cheshire, UK	2016	3	0	
	<i>Ecdyonurus venosus</i>	Cheshire, UK	2016	6	0	
	<i>Leptophlebia vespertina</i>	Hampshire, UK	2016	1	0	
	<i>Paraleptophlebia submarginata</i>	Stirling, Scotland, UK	2017	3	0	
	<i>Rhithrogena semicolorata</i>	Stirling, Scotland, UK	2017	3	0	
Trichoptera	<i>Hydropsyche</i> sp.	Stirling, Scotland, UK	2017	3	0	
	<i>Polycentropus flavomaculatus</i>	Cheshire, UK	2017	3	0	
	<b><i>Rhyacophila dorsalis</i></b>	<b>Stirling, Scotland, UK</b>	<b>2017</b>	<b>3</b>	<b>2</b>	
Plecoptera	<i>Amphinemura sulcicollis</i>	Stirling, Scotland, UK	2017	3	0	
	<i>Dinocras cephalotes</i>	Stirling, Scotland, UK	2017	3	0	
	<i>Isoperla grammatica</i>	Stirling, Scotland, UK	2017	3	0	
	<i>Perla bipunctata</i>	Stirling, Scotland, UK	2017	3	0	
Hemiptera	<i>Corixa punctata</i>	Cheshire, UK	2016	1	0	
	<i>Gerris</i> sp.	Montferrier sur Lez, France	2006	12	0	
	<i>Gerris thoracicus</i>	Cheshire, UK	2016	1	0	
	<i>Hydrometra stagnorum</i>	Montferrier sur Lez, France	2006	20	0	
	<i>Nepa cinerea</i>	Montferrier sur Lez, France	2006	3	0	
	<i>Notonecta glauca</i>	Cheshire, UK	2016	2	0	
	<i>Plea minutissima</i>	Notre Dame de Londres, France	2006	8	0	
	<i>Sigara lateralis</i>	Notre Dame de Londres, France	2006	6	0	
	<b><i>Sigara striata</i></b>	<b>Cheshire, UK</b>	<b>2006</b>	<b>2</b>	<b>1</b>	
Diptera	<i>Aedes</i> sp.	Cheshire, UK	2017	8	0	
	<i>Aedes albopictus</i>	Roma, Italy	2005	20	0	
	<b><i>Anopheles plumbeus</i></b>	<b>Chester Zoo, UK</b>	<b>2018</b>	<b>2</b>	<b>2</b>	
	<b>Chironomidae</b> sp.	<b>Cheshire, UK</b>	<b>2016</b>	<b>4</b>	<b>1</b>	
	<i>Chironomus acidophilus</i>	Cheshire, UK	2017	1	0	
	<i>Chironomus plumosus</i>	Notre Dame de Londres, France	2006	20	0	
	<i>Chironomus</i> sp.	Cheshire, UK	2016	4	0	
	<i>Culex pipiens</i> (ssp. quinquefasciatus)	Puerto Viejo de Talamanca, Costa Rica	2006	20	0	
	<i>Culex pipiens</i>	St Nazaire de Pézan, France	2006	20	0	
	<i>Eristalinus</i> sp.	Cheshire, UK	2016	3	0	
	<i>Eristalis tenax</i>	Montpellier (grotte du zoo), France	2002	7	0	
		<b><i>Glyptotendipes</i> sp.</b>	<b>Cheshire, UK</b>	<b>2016</b>	<b>1</b>	<b>1</b>
		<b><i>Hilara interstincta</i></b>	<b>Cheshire, UK</b>	<b>2017</b>	<b>3</b>	<b>1</b>
		<b><i>Simulium aureum</i></b>	<b>Hampshire, UK</b>	<b>2017</b>	<b>1</b>	<b>1</b>
		<i>Simulium ornatum</i>	N/A	2003	12	0
		<i>Tipula</i> sp.	UK	2006	10	0
	<i>Tipula oleracea</i>	UK	2006	13	0	
	<b><i>Zavrelimyia</i> sp.</b>	<b>Northumberland, UK</b>	<b>2017</b>	<b>1</b>	<b>1</b>	
Coleoptera	<i>Agabus bipustulatus</i>	Cheshire, UK	2017	3	0	
	<i>Guignotus pusillus</i>	Notre Dame de Londres, France	2006	12	0	
	Unknown sp.1	Cheshire, UK	2017	2	0	
	Unknown sp.2	Cheshire, UK	2017	3	0	
Acarina	Unknown sp.	Cheshire, UK	2017	3	0	
Isopoda	<i>Asellus aquaticus</i>	Cheshire, UK	2016	3	0	
Amphipoda	<i>Gammarus pulex</i>	Stirling, Scotland, UK	2017	3	0	
	<i>Crangonyx pseudogracilis</i>	Cheshire, UK	2016	6	0	
Gastropoda	<i>Radix balthica</i>	Cheshire, UK	2016	3	0	
	<i>Planorbis</i> sp.	Cheshire, UK	2016	3	0	
	<b><i>Galba truncatula</i></b>	<b>Cheshire, UK</b>	<b>2017</b>	<b>20</b>	<b>3</b>	
Hirudinea	<i>Erpobdella octoculata</i>	Cheshire, UK	2016	2	0	
	<i>Hemiclepsis marginata</i>	Cheshire, UK	2017	1	0	
Tricladida	Unknown sp.	Cheshire, UK	2016	1	0	

624 A species was deemed positive through PCR and designated to *Rickettsia* group after Sanger  
 625 sequencing and phylogenetic placement. All strains belong to the Torix group.

626  
 627  
 628  
 629

**Table 1.2.** Targeted *Rickettsia* screen of terrestrial invertebrates.

Terrestrial Invertebrate group	Species	Location	Year	Number tested	Number positive
Araneae	<i>Agelenopsis aperta</i>	Tennessee, USA	N/A	12	0
	<i>Allopecosa pulverulenta</i>	Berne, Germany	N/A	16	0
	<b><i>Amaurobius fenestralis</i></b>	<b>Montpellier, France</b>	<b>2006</b>	<b>16</b>	<b>1</b>
	<i>Araneus diadematus</i>	Beerse, Belgium	N/A	19	0
	<i>Araneus diadematus</i>	Greater London, UK	N/A	8	0
	<i>Argiope bruennichi</i>	Hamburg, Germany	N/A	7	0
	<i>Argiope lobata</i>	Spain	N/A	7	0
	<i>Argiope lobata</i>	Israel	N/A	4	0
	<i>Cyclosa conica</i>	Brandenburg, Germany	N/A	11	0
	<i>Dysdera crocata</i>	Montpellier, France	2006	2	0
	<i>Enoplognatha ovata</i>	Greater London, UK	N/A	20	0
	<i>Erigone atra</i>	Cheshire, UK	2017	1	0
	<i>Evarcha falcata</i>	Beerse, Belgium	N/A	5	0
	<i>Holochnemus pluchei</i>	Montpellier, France	2006	7	0
	<b><i>Hylyphantes graminicola</i></b>	<b>Cheshire, UK</b>	<b>2017</b>	<b>1</b>	<b>1</b>
	<i>Larinioides cornutus</i>	Greater London, UK	N/A	6	0
	<i>Larinioides scolopetarius</i>	Hamburg, Germany	N/A	17	0
	<b><i>Linyphia triangularis</i></b>	<b>Berlin, Germany</b>	<b>N/A</b>	<b>9</b>	<b>9</b>
	<i>Linyphia triangularis</i>	Greater London, UK	N/A	6	0
	<i>Lycosa</i> sp.	Cheshire, UK	2017	2	0
	<i>Metellina mengei</i>	Greater London, UK	N/A	13	0
	<i>Metellina segmentata</i>	Brandenburg, Germany	N/A	9	0
	<i>Neriere clathrata</i>	Beerse, Belgium	N/A	13	0
	<i>Neriere peltata</i>	Cheshire, UK	2017	1	0
	<i>Pachygnatha degeeri</i>	Berne, Germany	N/A	11	0
	<i>Pachygnatha listeri</i>	Beerse, Belgium	N/A	17	0
	<b><i>Pardosa lugubris</i></b>	<b>Darmstadt, Germany</b>	<b>N/A</b>	<b>20</b>	<b>1</b>
	<i>Pardosa pullata</i>	Brandenburg, Germany	N/A	20	0
	<i>Pardosa purbeckensis</i>	Belgium	N/A	19	0
	<b><i>Pholcus phalangioides</i></b>	<b>Berlin, Germany</b>	<b>N/A</b>	<b>20</b>	<b>17</b>
<b><i>Pisaura mirabilis</i></b>	<b>Greater London, UK</b>	<b>N/A</b>	<b>12</b>	<b>1</b>	
<i>Tetragnatha montana</i>	Greater London, UK	N/A	20	0	
<i>Tetragnatha</i> sp.	Hampshire, UK	2017	3	0	
Unknown sp.	Cheshire, UK	2017	2	0	
<i>Xysticus cristatus</i>	Cambridgeshire, UK	N/A	16	0	
Opiliones	<i>Leiobunum rotundum</i>	Feurs, France	2006	6	0
Ixodida	<i>Ixodes uriae</i>	Hornøya, Norway	2005	19	0
	<i>Rhipicephalus microplus</i>	New Caledonia, France	2003	1	0
Scorpiones	<i>Euscorpis flavicauda</i>	St Nazaire de Pézan, France	2006	1	0
Diplopoda	<i>Ommatoiulus</i> sp.	Cheshire, UK	2016	1	0
Neuroptera	Unknown sp.	Cheshire, UK	2017	1	0
Mecoptera	<i>Panorpa</i> sp.	Cheshire, UK	2017	2	0
Orthoptera	<i>Calliptamus italicus</i>	Notre Dame de Londres, France	2016	18	0
	<i>Chorthippus brunneus</i>	Uk	2006	20	0
	<i>Grylломорpha dalmatina</i>	Montpellier, France	2006	2	0

Blattaria	<i>Loboptera decipiens</i>	Montpellier, France	2006	17	0	
Mantodae	<i>Iris oratoria</i>	St Nazaire de Pézan, France	2006	6	0	
	<i>Mantis religiosa</i>	Feurs, France	2006	3	0	
Dermaptera	<i>Forficula Auricularia</i>	Feurs, France	2006	9	0	
Hemiptera	<i>Aphis fabae</i>	Montpellier, France	2006	12	0	
	<i>Aphis nerii</i>	Montpellier, France	2006	8	0	
	<i>Baizongia pistaciae</i>	Viols le Fort, France	2006	12	0	
	<i>Cicadella viridis</i>	L'Olme, France	2006	16	0	
	<b><i>Cimex lectularius</i></b>	<b>Yorkshire, UK</b>	<b>2008</b>	<b>12</b>	<b>12</b>	
	<i>Elasmucha grisea</i>	Greater London, UK	2006	16	0	
	<i>Graphosoma italicum</i>	Montpellier, France	2006	12	0	
	<i>Lygaeus equestris</i>	Montpellier, France	2006	12	0	
	<i>Notostira elongata</i>	L'Olme, France	2006	11	0	
	<i>Pyrrhocoris apterus</i>	Montpellier, France	2006	11	0	
	<i>Rhyparochromus vulgaris</i>	Castelnaudary, France	2006	20	0	
Coleoptera	<i>Anaspis frontalis</i>	Mont Barri, France	2004	12	0	
	<i>Anthaxia nitidula</i>	Mont Barri, France	2004	20	0	
	<i>Anthaxia</i> sp.	Mont Barri, France	2004	16	0	
	<i>Calvia 14-guttata</i>	Greater London, UK	2006	6	0	
	<i>Capnodis tenebrionis</i>	Montpellier, France	2006	1	0	
	<i>Cetonia aurata</i>	Feurs, France	2006	3	0	
	<i>Cetonia aurata</i>	Mont Barri, France	2004	12	0	
	<i>Chrysolina varians</i>	Mont Barri, France	2004	18	0	
	<i>Clytus arietis</i>	Mont Barri, France	2004	20	0	
	<i>Dermestes</i> sp.	Mont Barri, France	2004	20	0	
	<i>Dermestes tessellatocollis</i>	Cheshire, UK	2016	2	0	
	<i>Gastrophysa</i> sp.	Greater London, UK	2006	20	0	
	<i>Geotrupes stercorarius</i>	Mont Barri, France	2004	3	0	
	<i>Larinus scolymi</i>	Aldira de Irmeros, Spain	2005	12	0	
	<i>Leptinotarsa decemlineata</i>	Feurs, France	2006	10	0	
	<i>Mordellistena</i> sp.	Mont Barri, France	2004	10	0	
	<i>Oedemera</i> sp.	Mont Barri, France	2004	20	0	
		<i>Oncocerna</i> sp.	Mont Barri, France	2004	20	0
	<b><i>Phyllobius argentatus</i></b>	<b>Mont Barri, France</b>	<b>2004</b>	<b>15</b>	<b>4†</b>	
	<i>Pseudovadonia livida</i>	Mont Barri, France	2004	19	0	
	<i>Stenopterus</i> sp.	Mont Barri, France	2004	20	0	
Diptera	<i>Braula coeca</i>	Ouessant, France	2002	4	0	
	<i>Chorisops tunisiae</i>	Montpellier, France	2003	8	0	
	<i>Delia antiqua</i>	N/A	N/A	11	0	
	<i>Delia platura</i>	N/A	N/A	11	0	
	<i>Delia radicum</i>	N/A	N/A	10	0	
	<i>Gasterophilus intestinalis</i>	France	N/A	10	0	
	<i>Hippobosca equina</i>	Restinclières, France	2006	15	0	
	<i>Lonchoptera lutea</i>	Cheshire, UK	2017	3	0	
		<i>Medetera petrophila</i>	St Bauzille de Putois, France	2003	12	0
		<i>Musca domestica</i>	L'Olme, France	2006	20	0
		<i>Musca vitripennis</i>	Notre Dame de Londres, France	2003	8	0
		<i>Neomyia cornicina</i>	Notre Dame de Londres, France	2003	8	0
		<i>Protocalliphora</i> sp.	Corse, France	2003	2	0
		<b><i>Protocalliphora azurea</i></b>	<b>Montpellier, France</b>	<b>2005</b>	<b>12</b>	<b>12</b>
		<i>Psila rosae</i>	N/A	N/A	11	0
	<i>Stomoxys calcitrans</i>	Le Malzieu, France	2001	11	0	
Lepidoptera	<i>Chilo phragmitellus</i>	Feurs, France	2006	10	0	
	<i>Euplagia quadripunctaria</i>	Feurs, France	2006	2	0	
	<i>Pieris brassicae</i>	Feurs, France	2006	7	0	
	<i>Plodia interpunctella</i>	Montpellier, France	2006	12	0	



	<i>Thymelicus lineola</i>	Greater London, UK	2006	15	0
	<i>Thymelicus sylvestris</i>	Greater London, UK	2006	2	0
	<i>Triodia sylvina</i>	Montpellier, France	2006	4	0
Hymenoptera	<i>Amblyteles armatorius</i>	St Nazaire de Pézan, France	2006	1	0
	<i>Amegilla albigena</i>	St Nazaire de Pézan, France	2006	13	0
	<i>Amegilla ochroleuca</i>	St Nazaire de Pézan, France	2006	3	0
	<i>Anthidium florentinum</i>	St Nazaire de Pézan, France	2006	6	0
	<i>Apis mellifera</i>	UK	2006	9	0
	<i>Bombus terrestris</i>	North West, Switzerland	2006	20	0
	<i>Diplolepis rosae</i>	L'Olme, France	2006	2	0
	<i>Formica lugubris</i>	UK	2006	10	0
	<b><i>Pachycrepoideus sp.</i></b>	<b>UK</b>	<b>N/A</b>	<b>94</b>	<b>6‡</b>
	<i>Polistes dominulus</i>	St Nazaire de Pézan, France	2006	4	0
	<i>Polistes nimpha</i>	St Nazaire de Pézan, France	2006	19	0
<i>Sceliphron caementarium</i>	St Nazaire de Pézan, France	2006	3	0	

630

631 A species was deemed positive through PCR and designated to *Rickettsia* group after Sanger

632 sequencing and phylogenetic placement. All strains belong to the Torix group except

633 †=Rhyzobius and ‡=Belli.

634

635 **Table 2.** Torix *Rickettsia* hosts known to date alongside screening method.

636

Order	Host	Screening method	Reference
Amphipoda	<b><i>Paracalliope fluviatilis</i></b> <b>(Paracalliopiidae)</b>	<b>GenBank search</b>	<b>This study</b>
	<i>Paraleptamphopus sp.</i> (Paraleptamphopidae)	Barcoding	[33]
	Senticaudata sp.	Barcoding	[33]
Araneae	<b><i>Amaurobius fenestralis</i></b> <b>(Amaurobiidae)</b>	<b>Targeted PCR</b>	<b>This study</b>
	<i>Amaurobioides africana</i> (Anyphaenidae)	Barcoding	[32]
	<i>Araneus diadematus</i> (Araneidae)	Targeted PCR	[43]
	<i>Dysdera microdonta</i> (Dysderidae)	Barcoding	[31]
	<i>Linyphiidae spp.</i>	Targeted PCR	[43]

	<i>Linyphia triangularis</i> (Linyphiidae)	Targeted PCR	This study
	<i>Pardosa lugubris</i> (Lycosidae)	Targeted PCR	This study
	<i>Pholcus phalangioides</i> (Pholcidae)	Targeted PCR	This study
	<i>Pisaura mirabilis</i> (Pisauridae)	Targeted PCR	This study
	<i>Metellina mendei</i> (Tetragnathidae)	Targeted PCR	[43]
Coleoptera	<i>Deronectes</i> spp. (Dytiscidae)	Targeted PCR, FISH and TEM	[24]
	<i>Dytiscidae</i> sp.	Barcoding	This study
	<i>Stegobium paniceum</i> (Ptinidae)	Non-targeted (16S) PCR	[86]
	<i>Prionocyphon limbatus</i> (Scirtidae)	Barcoding	This study
	<i>Labidopullus appendiculatus</i> (Staphylinidae)	SRA search	This study
	<i>Platyusa sonomae</i> (Staphylinidae)	SRA search	This study
	<i>Pseudomimeceton antennatum</i> (Staphylinidae)	SRA search	This study
	<i>Staphylinidae</i> sp.	Barcoding	This study
	<i>Pimelia</i> sp. (Tenebrionidae)	GenBank search	This study
Dermaptera	<i>Forficula</i> sp. (Forficulidae)	GenBank search	This study
	unknown sp.	Barcoding	This study
Diplopoda	<i>Polydesmus complanatus</i> (Polydesmidae)	Targeted PCR	[87]
	unknown sp.	Barcoding	This study
	<i>Protocalliphora azurea</i> (Calliphoridae)	Targeted PCR	This study
	<i>Cecidomyiidae</i> sp.	Barcoding	This study
	<i>Chaoborus trivittatus</i> (Chaoboridae)	SRA search	This study
	<i>Mochlonyx cinctipes</i> (Chaoboridae)	SRA search	This study

Diptera	<b><i>Glyptotendipes</i> sp. (Chironomidae)</b>	Targeted PCR	This study
	<b><i>Zavrelimyia</i> sp. (Chironomidae)</b>	Targeted PCR	This study
	<i>Culicoides</i> spp. (Ceratopogonidae)	Targeted PCR and FISH	[27]
	<b><i>Anopheles plumbeus</i> (Culicidae)</b>	Targeted PCR	This study
	<i>Dolichopodidae</i> spp.	Targeted PCR	[44]
	<i>Empididae</i> spp.	Targeted PCR	[44]
	<i>Limonia chorea</i> (Limoniidae)	N/A	Unpublished (AF322443)
	<b><i>Boletina villosa</i> (Mycetophilidae)</b>	Barcoding	This study
	<b><i>Gnoriste bilineata</i> (Mycetophilidae)</b>	SRA search	This study
	<b><i>Mycetophila lunata</i> (Mycetophilidae)</b>	GenBank search	This study
	<b><i>Psilidae</i> sp.</b>	Barcoding	This study
	<i>Lutzomyia apache</i> (Psychodidae)	Targeted PCR	[61]
	<i>Phlebotomus chinensis</i> (Psychodidae)	Non-targeted (16S) PCR	[60]
	<b><i>Sciaridae</i> sp.</b>	Barcoding	This study
	<b><i>Pherbellia tenuipes</i> (Sciomyzidae)</b>	Barcoding	This study
	<b><i>Simulium aureum</i> (Simuliidae)</b>	Targeted PCR	This study
<b><i>Tabanidae</i> sp.</b>	Barcoding	This study	
Gastropoda	<b><i>Galba truncatula</i> (Lymnaeidae)</b>	Targeted PCR	This study
Haplotaxida	<i>Mesenchytraeus solifugus</i> (Enchytraediae)	Non-targeted (16S) PCR	[88]
	<i>Bemisia tabaci</i> (Aleyrodidae)	Targeted PCR and FISH	[51]
	<i>Nephotettix cincticeps</i> (Cicadellidae)	Targeted PCR, FISH and TEM	[89]
	<i>Platypleura kaempferi</i> (Cicadidae)	Non-targeted (16S) PCR	[90]
	<b><i>Cimex lectularius</i> (Cimicidae)</b>	Targeted PCR	This study/[65]
	<b><i>Sigara striata</i> (Corixidae)</b>	Targeted PCR	This study

Hemiptera	<b><i>Metcalfa pruinosa</i></b> (Flatidae)	GenBank search	This study
	<b><i>Flavina</i> sp.</b> (Issidae)	GenBank search	This study
	<i>Centrotus cornutus</i> (Membracidae)	Non-targeted (16S) PCR and TEM	[91]
	<i>Gargara genistae</i> (Membracidae)	Non-targeted (16S) PCR and TEM	[91]
	<i>Macrolophus pygmaeus</i> (Miridae)	Non-targeted (16S) PCR and FISH	[45]
	<b><i>Cacopsylla melanoneura</i></b> (Psyllidae)	Barcoding	This study
	<b><i>Chamaepsylla hartigii</i></b> (Psyllidae)	Barcoding	This study
	<b><i>Ricaniidae</i> sp.</b>	Barcoding	This study
Hirudinea	<i>Hemiclepsis</i> spp. (Glossiphoniidae)	Targeted PCR and TEM	[23]
	<i>Torix</i> spp. (Glossiphoniidae)	Targeted PCR and TEM	[23]
Hymenoptera	<i>Asobara tabida</i> (Braconidae)	Non-targeted (16S) PCR	[92]
	<b><i>Ceraphronidae</i> sp.</b>	Barcoding	This study
	<b><i>Diapriidae</i> sp.</b>	Barcoding	This study
	<b><i>Eucharitidae</i> sp.</b>	GenBank search	This study
	<i>Quadrastichus mendeli</i> (Eulophidae)	Non-targeted (16S) PCR and FISH	[93]
	<b><i>Formicidae</i> sp.</b>	GenBank search	This study
	<i>Atta colombica</i> (Formicidae)	Non-targeted (16S) PCR	Unpublished (LN570502)
	<b><i>Megaspilidae</i> sp.</b>	Barcoding	This study
	<b><i>Mymaridae</i> sp.</b>	Barcoding	This study
<b><i>Platygastridae</i> sp.</b>	Barcoding	This study	
Ixodida	<i>Argas japonica</i> (Argasidae)	Non-targeted (16S) PCR	[64]
	<i>Ixodes ricinus</i> (Ixodidae)	Targeted PCR	[63]
Megaloptera	<i>Sialis lutaria</i> (Sialidae)	Targeted PCR	[94]
Neuroptera	<i>Chrysotropia ciliata</i> (Chrysopidae)	Targeted PCR	[94]

Nucleariida	<i>Nuclearia pattersoni</i> (Nucleariidae)	Non-targeted (16S) PCR	[25]
	<i>Pompholyxophrys punicea</i> (Pompholyxophryidae)	Single cell sequencing	[26]
Odonata	<b><i>Calopteryx maculata</i></b> <b>(Calopterygidae)</b>	<b>GenBank search</b>	<b>This study</b>
	<i>Coenagrionidae</i> spp.	Targeted PCR and FISH	[28]
	<i>Sympetrum fonscolombii</i> (Libellulidae)	Targeted PCR	[28]
	<i>Polythoridae</i> spp.	Targeted PCR	[28]
	<i>Neoneura sylvatica</i> (Protoneuridae)	Targeted PCR	[28]
Psocoptera	<b><i>Myopsocidae</i> sp.</b>	<b>Barcoding</b>	<b>This study</b>
	<b><i>Philotarsus californicus</i></b> <b>(Philotarsidae)</b>	<b>Barcoding</b>	<b>This study</b>
	<i>Cerobasis guestfalica</i> (Trogidae)	Targeted PCR and FISH	[95]
Siphonaptera	<i>Nosopsyllus fasciatus</i> (Ceratophyllidae)	Targeted PCR	[62]
Trichoptera	<b><i>Lepidostoma hoodi</i></b> <b>(Lepidostomatidae)</b>	<b>Barcoding</b>	<b>This study</b>
	<b><i>Rhyacophila dorsalis</i></b> <b>(Rhyacophilidae)</b>	<b>Targeted PCR</b>	<b>This study</b>
	<b><i>Sericostoma</i> sp.</b> <b>(Sericostomatidae)</b>	<b>SRA search</b>	<b>This study</b>
Trombidiformes	<b><i>Calyptostomatidae</i> sp.</b>	<b>Barcoding</b>	<b>This study</b>

637

638 Bold entries indicate hosts identified in this study. FISH=fluorescence *in-situ* hybridisation;

639 TEM=transmission electron microscopy; SRA=sequence read archive. Accession numbers for

640 *Rickettsia* sequences from newly detected hosts can be found in Additional files 8 and 10.

641

## 642 **Availability of Supporting Data and Materials**

643 The data sets supporting the findings of this study are openly available in:

644 The Barcode of Life Data System (BOLD) repository [37] and the Figshare repository [36][39].

645 Alignments and trees are also available from the *GigaScience* GigaDB repository [96].

646 For DNA sequences, accessions are: Bioproject number PRJEB38316; LR798809-LR800243;  
647 LR812141-LR812260; LR812269-LR812283; LR812678; LR813674-LR813676; LR813730.

648

## 649 **Declarations**

### 650 **List of Abbreviations**

651 BOLD = Barcode of Life Data System

652 COI = cytochrome c oxidase I

653 FISH = fluorescence *in-situ* hybridisation

654 SRA = Sequence Read Archive

655

### 656 **Ethics Approval**

657 Not applicable.

658

### 659 **Consent for Publication**

660 Not applicable.

661

### 662 **Competing Interests**

663 The authors declare that they have no competing interests.

664

### 665 **Funding**

666 This work was supported by: a BBSRC Doctoral Training Partnership studentship

667 (BB/M011186/1) awarded to JP; a Development and Promotion of Science and Technology

668 Talents Project (DPST) of the Institute for the Promotion of Teaching Science and Technology,

669 Thailand to PT and a Harry Smith Vacation studentship (Microbiology society) and a NERC  
670 ACCE DTP studentship (NE/L002450/1) to HRD.

671

## 672 **Acknowledgments**

673 We would like to thank Dr. Michael Gerth for kindly providing comments on the manuscript.

674

## 675 **Author contributions**

676 JP, GDDH, MB and MAS: conception and design of the study. MAS, EVZ, SR and JRD:

677 assembling BOLD datasets and providing DNA extracts for laboratory experiments. Field and

678 laboratory work: JP, CRM and PT. SRA work: HRD and SS. Analyses and interpretation of the

679 data, drafting of the manuscript: JP, PT, HRD, GDDH, MB and SS. All authors assisted in

680 critical revision of the manuscript.

681

## 682 **References**

683 1. McFall-Ngai M, Hadfield MG, Bosch TCG, Carey H V., Domazet-Lošo T, Douglas AE, et al.

684 Animals in a bacterial world, a new imperative for the life sciences. *Proc Natl Acad Sci.*

685 2013;110:3229–36.

686 2. Hurst GDD. Extended genomes: symbiosis and evolution. *Interface Focus.* 2017;7:20170001.

687 3. Łukasik P, Guo H, van Asch M, Ferrari J, Godfray HCJ. Protection against a fungal pathogen

688 conferred by the aphid facultative endosymbionts *Rickettsia* and *Spiroplasma* is expressed in

689 multiple host genotypes and species and is not influenced by co-infection with another

690 symbiont. *J Evol Biol.* 2013;26:2654–61.

- 691 4. Teixeira L, Ferreira A, Ashburner M. The bacterial symbiont Wolbachia induces resistance  
692 to RNA viral infections in *Drosophila melanogaster*. *PLoS Biol.* 2008;6:2753–63.
- 693 5. Rio RVM, Attardo GM, Weiss BL. Grandeur Alliances: Symbiont metabolic integration and  
694 obligate arthropod hematophagy. *Trends Parasitol.* 2016;32:739–49.
- 695 6. Douglas AE. The microbial dimension in insect nutritional ecology. *Funct Ecol.* 2009;23:38–  
696 47.
- 697 7. Hurst GDD, Frost CL. Reproductive parasitism: Maternally inherited symbionts in a  
698 biparental world. *Cold Spring Harb Perspect Biol.* 2015;7:a017699.
- 699 8. Munson MA, Baumann P, Kinsey MG. *Buchnera* gen. nov. and *Buchnera aphidicola* sp. nov.,  
700 a taxon consisting of the mycetocyte-associated, primary endosymbionts of aphids. *Int J Syst*  
701 *Bacteriol.* 1991;41:566–8.
- 702 9. Zug R, Hammerstein P. Still a host of hosts for Wolbachia: Analysis of recent data suggests  
703 that 40% of terrestrial arthropod species are infected. *PLoS One.* 2012;7:e38544.
- 704 10. Siozios S, Gerth M, Griffin JS, Hurst GDD. Symbiosis: Wolbachia host shifts in the fast lane.  
705 *Curr Biol.* 2018;28:R269–71.
- 706 11. Duron O, Bouchon D, Boutin S, Bellamy L, Zhou L, Engelstädter J, et al. The diversity of  
707 reproductive parasites among arthropods: Wolbachia do not walk alone. *BMC Biol.* 2008;6:27.
- 708 12. Weinert LA, Araujo-Jnr E V, Ahmed MZ, Welch JJ. The incidence of bacterial endosymbionts  
709 in terrestrial arthropods. *Proc R Soc B.* 2015;282:20150249.
- 710 13. Weinert LA, Werren JH, Aebi A, Stone GN, Jiggins FM. Evolution and diversity of *Rickettsia*  
711 bacteria. *BMC Biol.* 2009;7:6.
- 712 14. Perlman SJ, Hunter MS, Zchori-Fein E. The emerging diversity of *Rickettsia*. *Proc R Soc B.*  
713 2006;273:2097–106.



- 714 15. Ricketts HT. A micro-organism which apparently has a specific relationship to Rocky  
715 Mountain spotted fever. *J Am Med Assoc.* 1909;52:379–80.
- 716 16. da Rocha-Lima H. Zur Aetiologie des Fleckfiebers. *Dtsch Medizinische Wochenschrift.*  
717 1916;53:567–9.
- 718 17. Werren JH, Hurst GD, Zhang W, Breeuwer JA, Stouthamer R, Majerus ME. Rickettsial  
719 relative associated with male killing in the ladybird beetle (*Adalia bipunctata*). *J Bacteriol.*  
720 1994;176:388–94.
- 721 18. Chen D-Q, Campbell BC, Purcell AH. A new *Rickettsia* from a herbivorous insect, the pea  
722 aphid *Acyrtosiphon pisum* (Harris). *Curr Microbiol.* 1996;33:123–8.
- 723 19. Hurst GDD, Walker LE, Majerus MEN. Bacterial infections of hemocytes associated with  
724 the maternally inherited male-killing trait in British populations of the two spot ladybird,  
725 *Adalia bipunctata*. *J Invertebr Pathol.* 1996;68:286–92.
- 726 20. Hendry TA, Hunter MS, Baltrus DA. The facultative symbiont *Rickettsia* protects an invasive  
727 whitefly against entomopathogenic *Pseudomonas syringae* strains. *Appl Environ Microbiol.*  
728 2014;80:7161–8.
- 729 21. Himler AG, Adachi-Hagimori T, Bergen JE, Kozuch A, Kelly SE, Tabashnik BE, et al. Rapid  
730 spread of a bacterial symbiont in an invasive whitefly is driven by fitness benefits and female  
731 bias. *Science.* 2011;332:254–6.
- 732 22. Kikuchi Y, Fukatsu T. *Rickettsia* infection in natural leech populations. *Microb Ecol.*  
733 2005;49:265–71.
- 734 23. Kikuchi Y, Sameshima S, Kitade O, Kojima J, Fukatsu T. Novel clade of *Rickettsia* spp. from  
735 leeches. *Appl Environ Microbiol.* 2002;68:999–1004.

- 736 24. Kűchler SM, Kehl S, Dettner K. Characterization and localization of *Rickettsia* sp. in water  
737 beetles of genus *Deronectes* (Coleoptera: Dytiscidae). *FEMS Microbiol Ecol.* 2009;68:201–11.
- 738 25. Dyková I, Veverková M, Fiala I, Macháčková B, Pecková H. *Nuclearia pattersoni* sp. n.  
739 (Filosea), a new species of amphizoic amoeba isolated from gills of roach (*Rutilus rutilus*), and  
740 its Rickettsial endosymbiont. *Folia Parasitol (Praha).* 2003;50:161–70.
- 741 26. Galindo LJ, Torruella G, Moreira D, Eglit Y, Simpson AGB, Völcker E, et al. Combined  
742 cultivation and single-cell approaches to the phylogenomics of nucleariid amoebae, close  
743 relatives of fungi. *Philos Trans R Soc B Biol Sci.* 2019;374:20190094.
- 744 27. Pilgrim J, Ander M, Garros C, Baylis M, Hurst GDD, Siozios S. Torix group *Rickettsia* are  
745 widespread in *Culicoides* biting midges (Diptera: Ceratopogonidae), reach high frequency and  
746 carry unique genomic features. *Environ Microbiol.* 2017;19:4238–55.
- 747 28. Thongprem P, Davison HR, Thompson DJ, Lorenzo-Carballa MO, Hurst GDD. Incidence and  
748 diversity of Torix *Rickettsia*–Odonata symbioses. *Microb Ecol.* 2020; DOI:10.1007/s00248-020-  
749 01568-9
- 750 29. Weinert LA. The diversity and phylogeny of *Rickettsia*. In: Morand S, Krasnov BR,  
751 Littlewood DTJ, editors. *Parasite diversity and diversification*. Cambridge: Cambridge  
752 University Press; 2015. p. 150–81.
- 753 30. Lagrue C, Joannes A, Poulin R, Blasco-Costa I. Genetic structure and host-parasite co-  
754 divergence: evidence for trait-specific local adaptation. *Biol J Linn Soc.* 2016;118:344–58.
- 755 31. Řezáč M, Gasparo F, Král J, Heneberg P. Integrative taxonomy and evolutionary history of  
756 a newly revealed spider *Dysdera ninnii* complex (Araneae: Dysderidae). *Zool J Linn Soc.*  
757 2014;172:451–74.

- 758 32. Ceccarelli FS, Haddad CR, Ramírez MJ. Endosymbiotic Rickettsiales (Alphaproteobacteria)  
759 from the spider genus Amaurobioides (Araneae: Anyphaenidae). J Arachnol. 2016;44:251–3.
- 760 33. Park E, Poulin R. Widespread *Rickettsia* in New Zealand amphipods and the use of  
761 blocking primers to rescue host COI sequences. Sci Rep. 2020;10:16842.
- 762 34. Smith MA, Bertrand C, Crosby K, Eveleigh ES, Fernandez-Triana J, Fisher BL, et al.  
763 Wolbachia and DNA barcoding insects: Patterns, potential, and problems. PLoS One.  
764 2012;7:e36514.
- 765 35. BOLD: Barcode of Life Data System. 2007. <https://www.boldsystems.org/>  
766 Accessed 2 January 2018.
- 767 36. Smith MA, Pilgrim J, Zakharov E V., Dewaard JR, Ratnasingham S. BOLD contaminant pool  
768 (3,817 specimens) data. Figshare. 2020; DOI:10.6084/m9.figshare.12801107
- 769 37. Smith MA, Pilgrim J, Zakharov E V., Dewaard JR, Ratnasingham S. BOLD non-contaminant  
770 pool (55,366 specimens) data. Barcode Of Life Data System. 2020; DOI:10.5883/DS-RICKET
- 771 38. Ratnasingham S, Hebert PDN. BOLD: The Barcode of Life Data System. Mol Ecol Notes.  
772 2007;7:355–64.
- 773 39. Davison HR, Siozios S. *Rickettsia* PhyloFlash and Kraken data from arthropod whole  
774 genome projects in the Sequence Read Archive. Figshare. 2020;  
775 DOI:10.6084/m9.figshare.12801140
- 776 40. Gruber-Vodicka HR, Seah BK, Pruesse E. phyloFlash: Rapid small-subunit rRNA profiling and  
777 targeted assembly from metagenomes. mSystems. 2020; DOI:10.1128/mSystems.00920-20
- 778 41. Hernández-Triana LM, Prosser SW, Rodríguez-Perez MA, Chaverri LG, Hebert PDN, Ryan  
779 Gregory T. Recovery of DNA barcodes from blackfly museum specimens (Diptera: Simuliidae)  
780 using primer sets that target a variety of sequence lengths. Mol Ecol Resour. 2014;14:508–18.

- 781 42. Wang Z, Wu M. An integrated phylogenomic approach toward pinpointing the origin of  
782 mitochondria. *Sci Rep.* 2015;5:7949.
- 783 43. Goodacre SL, Martin OY, Thomas CFG, Hewitt GM. *Wolbachia* and other endosymbiont  
784 infections in spiders. *Mol Ecol.* 2006;15:517–27.
- 785 44. Martin OY, Puniamoorthy N, Gubler A, Wimmer C, Bernasconi M V. Infections with  
786 *Wolbachia*, *Spiroplasma*, and *Rickettsia* in the Dolichopodidae and other Empidoidea. *Infect*  
787 *Genet Evol.* 2013;13:317–30.
- 788 45. Machtelinckx T, Van Leeuwen T, Van De Wiele T, Boon N, De Vos WH, Sanchez J-A, et al.  
789 Microbial community of predatory bugs of the genus *Macrolophus* (Hemiptera: Miridae). *BMC*  
790 *Microbiol.* 2012;12:S9.
- 791 46. Lanzoni O, Sabaneyeva E, Modeo L, Castelli M, Lebedeva N, Verni F, et al. Diversity and  
792 environmental distribution of the cosmopolitan endosymbiont “*Candidatus Megaira*”. *Sci Rep.*  
793 2019;9:1179.
- 794 47. Blaxter M. Symbiont genes in host genomes: Fragments with a future? *Cell Host Microbe.*  
795 2007;2:211-3.
- 796 48. Gehrler L, Vorburger C. Parasitoids as vectors of facultative bacterial endosymbionts in  
797 aphids. *Biol Lett.* 2012;8:613-5.
- 798 49. Le Clec’h W, Chevalier FD, Genty L, Bertaux J, Bouchon D, Sicard M. Cannibalism and  
799 predation as paths for horizontal passage of *Wolbachia* between terrestrial isopods.
- 800 50. Ramage T, Martins-Simoes P, Mialdea G, Allemand R, Duplouy A, Rousse P, et al. A DNA  
801 barcode-based survey of terrestrial arthropods in the Society Islands of French Polynesia: host  
802 diversity within the SymbioCode Project. *Eur J Taxon.* 2017;272.

- 803 51. Wang H, Lei T, Wang X, Maruthi MN, Zhu D, Cameron SL, et al. A newly recorded *Rickettsia*  
804 of the Torix group is a recent intruder and an endosymbiont in the whitefly *Bemisia tabaci*.  
805 *Environ Microbiol.* 2020;22:1207–21.
- 806 52. Caspi-Fluger A, Inbar M, Mozes-Daube N, Katzir N, Portnoy V, Belausov E, et al. Horizontal  
807 transmission of the insect symbiont *Rickettsia* is plant-mediated. *Proc R Soc B Biol Sci.*  
808 2012;279:1791–6.
- 809 53. Gonella E, Pajoro M, Marzorati M, Crotti E, Mandrioli M, Pontini M, et al. Plant-mediated  
810 interspecific horizontal transmission of an intracellular symbiont in insects. *Sci Rep.*  
811 2015;5:15811.
- 812 54. Li Y-H, Ahmed MZ, Li S-J, Lv N, Shi P-Q, Chen X-S, et al. Plant-mediated horizontal  
813 transmission of *Rickettsia* endosymbiont between different whitefly species. *FEMS Microbiol*  
814 *Ecol.* 2017;93.
- 815 55. Jaenike J, Polak M, Fiskin A, Helou M, Minhas M. Interspecific transmission of  
816 endosymbiotic *Spiroplasma* by mites. *Biol Lett.* 2007;3:23–5.
- 817 56. Morrow JL, Frommer M, Shearman DCA, Riegler M. Tropical tephritid fruit fly community  
818 with high incidence of shared *Wolbachia* strains as platform for horizontal transmission of  
819 endosymbionts. *Environ Microbiol.* 2014;16:3622–37.
- 820 57. Goodacre SL, Martin OY, Bonte D, Hutchings L, Woolley C, Ibrahim K, et al. Microbial  
821 modification of host long-distance dispersal capacity. *BMC Biol.* 2009;7:32.
- 822 58. Giorgini M, Bernardo U, Monti MM, Nappo AG, Gebiola M. *Rickettsia* symbionts cause  
823 parthenogenetic reproduction in the parasitoid wasp *Pnigalio soemius* (Hymenoptera:  
824 Eulophidae). *Appl Environ Microbiol.* 2010;76:2589–99.

- 825 59. Hagimori T, Abe Y, Date S, Miura K. The first finding of a *Rickettsia* bacterium associated  
826 with parthenogenesis induction among insects. *Curr Microbiol.* 2006;52:97–101.
- 827 60. Li K, Chen H, Jiang J, Li X, Xu J, Ma Y. Diversity of bacteriome associated with *Phlebotomus*  
828 *chinensis* (Diptera: Psychodidae) sand flies in two wild populations from China. *Sci Rep.*  
829 2016;6:36406.
- 830 61. Reeves WK, Kato CY, Gilchrist T. Pathogen screening and bionomics of *Lutzomyia apache*  
831 (Diptera: Psychodidae) in Wyoming, USA. *J Am Mosq Control Assoc.* 2008;24:444–7.
- 832 62. Song S, Chen C, Yang M, Zhao S, Wang B, Hornok S, et al. Diversity of *Rickettsia* species in  
833 border regions of northwestern China. *Parasit Vectors.* 2018;11:634.
- 834 63. Floris R, Yurtman AN, Margoni EF, Mignozzi K, Boemo B, Altobelli A, et al. Detection and  
835 identification of *Rickettsia* species in the Northeast of Italy. *Vector-Borne Zoonotic Dis.*  
836 2008;8:777–82.
- 837 64. Yan P, Qiu Z, Zhang T, Li Y, Wang W, Li M, et al. Microbial diversity in the tick *Argas*  
838 *japonicus* (Acari: Argasidae) with a focus on *Rickettsia* pathogens. *Med Vet Entomol.*  
839 2019;33:327–35.
- 840 65. Potts R, Molina I, Sheele JM, Pietri JE. Molecular detection of *Rickettsia* infection in field-  
841 collected bed bugs. *New Microbes New Infect.* 2020;34:100646.
- 842 66. Parola P, Paddock CD, Raoult D. Tick-borne Rickettsioses around the world: Emerging  
843 diseases challenging old concepts. *Clin Microbiol Rev.* 2005;18:719–56.
- 844 67. Hoffmann AA, Ross PA, Rašić G. Wolbachia strains for disease control: ecological and  
845 evolutionary considerations. *Evol Appl.* 2015;8:751–68.
- 846 68. Iturbe-Ormaetxe I, Walker T, O’Neill SL. Wolbachia and the biological control of mosquito-  
847 borne disease. *EMBO Rep.* 2011;12:508–18.

- 848 69. van den Hurk AF, Hall-Mendelin S, Pyke AT, Frentiu FD, McElroy K, Day A, et al. Impact of  
849 Wolbachia on infection with chikungunya and yellow fever viruses in the mosquito vector  
850 *Aedes aegypti*. PLoS Negl Trop Dis. 2012;6.
- 851 70. Kliot A, Cilia M, Czosnek H, Ghanim M. Implication of the bacterial endosymbiont *Rickettsia*  
852 spp. in interactions of the whitefly *Bemisia tabaci* with Tomato yellow leaf curl virus. J Virol.  
853 2014;88:5652–60.
- 854 71. Tedeschi R, Visentin C, Alam A, Bosco D. Epidemiology of apple proliferation (AP) in  
855 northwestern Italy: evaluation of the frequency of AP-positive psyllids in naturally infected  
856 populations of *Cacopsylla melanoneura* (Homoptera: Psyllidae). Ann Appl Biol. 2003;142:285-  
857 90.
- 858 72. Menzel P, Ng KL, Krogh A. Fast and sensitive taxonomic classification for metagenomics  
859 with Kaiju. Nat Commun. 2016;7:11257.
- 860 73. Katoh K, Standley DM. MAFFT Multiple sequence alignment software version 7:  
861 Improvements in performance and usability. Mol Biol Evol. 2013;30:772–80.
- 862 74. Castresana J. Selection of conserved blocks from multiple alignments for their use in  
863 phylogenetic analysis. Mol Biol Evol. 2000;17:540–52.
- 864 75. Wiens J. Missing data and the design of phylogenetic analyses. J. Biomed. Inform.  
865 2006;39:34-42.
- 866 76. Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermiin LS. ModelFinder: fast  
867 model selection for accurate phylogenetic estimates. Nat Methods. 2017;14:587–9.
- 868 77. Nguyen LT, Schmidt HA, Von Haeseler A, Minh BQ. IQ-TREE: A fast and effective stochastic  
869 algorithm for estimating maximum-likelihood phylogenies. Mol Biol Evol. 2015;32:268–74.

- 870 78. Hoang DT, Chernomor O, Haeseler A von, Minh BQ, Vinh LS. UFBoot2: Improving the  
871 ultrafast bootstrap approximation. *Mol Biol Evol.* 2017;35:518–22.
- 872 79. He Z, Zhang H, Gao S, Lercher MJ, Chen W-H, Hu S. Evolview v2: an online visualization and  
873 management tool for customized and annotated phylogenetic trees. *Nucleic Acids Res.*  
874 2016;44:W236–41.
- 875 80. Fredslund J, Schauer L, Madsen LH, Sandal N, Stougaard J. PriFi: using a multiple alignment  
876 of related sequences to find primers for amplification of homologs. *Nucleic Acids Res.*  
877 2005;33:W516–20.
- 878 81. Martin DP, Murrell B, Golden M, Khoosal A, Muhire B. RDP4: Detection and analysis of  
879 recombination patterns in virus genomes. *Virus Evol.* 2015;1:1–5.
- 880 82. Chernomor O, von Haeseler A, Minh BQ. Terrace aware data structure for phylogenomic  
881 inference from supermatrices. *Syst Biol.* 2016;65:997–1008.
- 882 83. Okonechnikov K, Golosova O, Fursov M. Unipro UGENE: a unified bioinformatics toolkit.  
883 *Bioinformatics.* 2012;28:1166–7.
- 884 84. Ewing B, Hillier L, Wendl MC, Green P. Base-calling of automated sequencer traces using  
885 Phred. I. Accuracy Assessment. *Genome Res.* 1998;8:175–85.
- 886 85. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. *Genome Biol.*  
887 2019;20:257.
- 888 86. Kölsch G, Synefiaridou D. Shared ancestry of symbionts? *Sagrinae* and *Donaciinae*  
889 (Coleoptera, Chrysomelidae) harbor similar bacteria. *Insects.* 2012;3:473–91.
- 890 87. Li K, Stanojević M, Stamenković G, Ilić B, Paunović M, Lu M, et al. Insight into diversity of  
891 bacteria belonging to the order Rickettsiales in 9 arthropods species collected in Serbia. *Sci*  
892 *Rep.* 2019;9:18680.



- 893 88. Murakami T, Segawa T, Bodington D, Dial R, Takeuchi N, Kohshima S, et al. Census of  
894 bacterial microbiota associated with the glacier ice worm *Mesenchytraeus solifugus*. FEMS  
895 Microbiol Ecol. 2015;91.
- 896 89. Noda H, Watanabe K, Kawai S, Yukuhiro F, Miyoshi T, Tomizawa M, et al. Bacteriome-  
897 associated endosymbionts of the green rice leafhopper *Nephotettix cincticeps* (Hemiptera:  
898 Cicadellidae). Appl Entomol Zool. 2012;47:217–25.
- 899 90. Zheng Z, Wang D, He H, Wei C. Bacterial diversity of bacteriomes and organs of  
900 reproductive, digestive and excretory systems in two cicada species (Hemiptera: Cicadidae).  
901 PLoS One. 2017;12:e0175903.
- 902 91. Kobiałka M, Michalik A, Świerczewski D, Szklarzewicz T. Complex symbiotic systems of two  
903 treehopper species: *Centrotus cornutus* (Linnaeus, 1758) and *Gargara genistae* (Fabricius,  
904 1775) (Hemiptera: Cicadomorpha: Membracoidea: Membracidae). Protoplasma.  
905 2020;257:819–31.
- 906 92. Zouache K, Voronin D, Tran-Van V, Mavingui P. Composition of bacterial communities  
907 associated with natural and laboratory populations of *Asobara tabida* infected with  
908 *Wolbachia*. Appl Environ Microbiol. 2009;75:3755–64.
- 909 93. Gualtieri L, Nugnes F, Nappo AG, Gebiola M, Bernardo U. Life inside a gall: closeness does  
910 not favour horizontal transmission of *Rickettsia* between a gall wasp and its parasitoid. FEMS  
911 Microbiol Ecol. 2017;93.
- 912 94. Gerth M, Wolf R, Bleidorn C, Richter J, Sontowski R, Unrein J, et al. Green lacewings  
913 (Neuroptera: Chrysopidae) are commonly associated with a diversity of Rickettsial  
914 endosymbionts. Zool Lett. 2017;3:12.

915 95. Perotti MA, Clarke HK, Turner BD, Braig HR. *Rickettsia* as obligate and mycetomic  
916 bacteria. FASEB J. 2006;20:2372–4.  
917 96. Pilgrim J; Thongprem P; Davison HR; Siozios S; Baylis M; Zakharov EV; Ratnasingham S;  
918 deWaard JR; Macadam CR; Smith MA; Hurst GDD (2021): Supporting data for "*Torix*  
919 *Rickettsia* are widespread in arthropods and reflect a neglected symbiosis" GigaScience  
920 Database. <http://dx.doi.org/10.5524/100873>

921

## 922 **Figure Legends**

923 **Figure 1.** Workflow of the BOLD project demonstrating the acquisition and fates of  
924 contaminant and non-contaminant *COI* barcoding sequences.

925

926 **Figure 2.** Cladogram of the maximum likelihood (ML) tree of 1,126 proteobacteria *COI* contaminants  
927 retrieved from a BOLD project incorporating 184,585 arthropod specimens. The tree is based on 561  
928 bp and is rooted with the free-living alphaproteobacteria *Pelagibacter ubique*. Parentheses indicate  
929 the number of BOLD contaminants present in each group. Tips are labelled by BOLD processing ID and  
930 host arthropod taxonomy. The Rickettsiales genera of *Anaplasma*, *Rickettsia* (collapsed node),  
931 *Orientia* and *Wolbachia* supergroups (A, B, E and F), as well as the Legionellales genera *Legionella* and  
932 *Rickettsiella*, are included as reference sequences (Accession numbers: Additional file 10).

933

934 **Figure 3.** Cladogram of a maximum likelihood (ML) tree of 753 *COI Rickettsia* contaminants retrieved  
935 from a BOLD project incorporating 184,585 arthropod specimens. The tree is based on 561 bp and  
936 is rooted by the *Rickettsia* endosymbiont of *Ichthyophthirius multifiliis* (Candidatus Megaira) using  
937 the TVM+F+I+G4 model. Parentheses indicate the number of BOLD contaminants present in *Torix*

938 and non-Torix *Rickettsia* groups. Tips are labelled by BOLD processing ID and host arthropod  
939 taxonomy. The *Rickettsia* groups: Spotted fever, Transitional, Belli, Typhus, Rhyzobius and Torix are  
940 included as references (Accession numbers: Additional file 10).

941  
942 **Figure 4.** Phylogram of the maximum likelihood (ML) tree of 99 *COI Rickettsia* contaminants (prefix  
943 “BIOUG”) used for further phylogenetic analysis and 53 Non-BOLD reference profiles (Accession  
944 numbers: Additional file 10). The tree is based on the concatenation of 4 loci; *16S rRNA*, *17KDa*, *gltA*  
945 and *COI* under a partition model, with profiles containing at least 3 out of 4 sites included in the tree  
946 (2,834 bp total) and is rooted by *Rickettsia* endosymbiont of *Ichthyophthirius multifiliis* (*Candidatus*  
947 Megaira). Tips are labelled by host arthropod taxonomy.

948  
949 **Figure 5.** *16S rRNA* and *gltA* concatenated maximum likelihood (ML) phylogram (1,834 bp total)  
950 including *Rickettsia* hosts from SRA (Triangles) and targeted screens (Stars). The TIM3+F+R2 (16S)  
951 and K3Pu+F+G4 (*gltA*) models were chosen as best fitting models. Rooting is with *Orientia*  
952 *tsutsugamushi*. Accession numbers found in Additional file 10.

953  
954 **Figure 6.** Phylogram of a maximum likelihood (ML) tree of *COI Rickettsia* contaminants (prefix  
955 “BIOUG”) giving a host barcode and 43 Non-BOLD reference profiles. The tree is based on 4 loci;  
956 *16S rRNA*, *17KDa*, *gltA* and *COI* under a partition model with profiles containing at least 2 out of  
957 4 sites included in the tree (2,781 bp total) and is rooted by the *Rickettsia* endosymbiont of  
958 *Ichthyophthirius multifiliis* (*Candidatus* Megaira). The habitats and lifestyles of the host are given  
959 to the right of the phylogeny. Accession numbers found in Additional file 10.

960

961 **Additional file information**

962

963 **Additional file 1.docx** Taxonomic classification of BOLD non-target *COI* sequences via Kaiju.

964

965 **Additional file 2.7z** Rectangular phylogram trees of cladograms from Figures 2 and 3.

966

967 **Additional file 3.docx** Primer pairs involved in the unintended amplification of 753 *Rickettsia*  
968 *COI* from BOLD project.

969

970 **Additional file 4.docx** Homology of *Rickettsia* groups and *Wolbachia* to the most common  
971 forward primers (C\_LepFolF and C\_LepFolR) attributed to bacterial *COI* amplification from  
972 arthropod DNA extracts.

973

974 **Additional file 5.xlsx** Re-barcoding status and nearest BLAST hit of mtDNA *COI* arthropod DNA  
975 extracts accessed for further analysis, along with the success of multilocus *Rickettsia* profiles  
976 with allocated *Rickettsia* group (based on phylogenetic analysis) and co-infection status.

977

978 **Additional file 6.docx** The barcoding success rate of taxa which gave at least one bacteria *COI*  
979 inadvertent amplification (N=51,475 accessible specimens) with an adjusted *Rickettsia*  
980 frequency based on an estimated total number of arthropods to account for inaccessible  
981 specimens (N=125,402).

982

983 **Additional file 7.docx** Fisher's Exact analyses for comparison of Torix *Rickettsia* infection in  
984 aquatic versus terrestrial insects.

985

986 **Additional file 8.docx** GenBank matches mistaken for true mtDNA barcodes and their  
987 homology to *Rickettsia COI* (Accessed 29<sup>th</sup> June 2020).

988

989 **Additional file 9.pdf** Phylogram of a maximum likelihood (ML) tree of *COI Rickettsia* found in the  
990 GenBank database erroneously identified as mtDNA barcodes based on 577 bp. The HKY+F+G4  
991 model was chosen as the best fitting model using Modelfinder with the Bayesian information  
992 criterion (BIC).

993

994 **Additional file 10.xlsx** Accession numbers used for phylogenetic analyses (Figures 2, 3, 4 ,5  
995 and 6). Accession numbers generated in this study are marked in BOLD.

996

997 **Additional file 11.docx** Mitochondrial *COI* and bacterial gene primers used for re-barcoding  
998 and multilocus phylogenetic analyses.

999

1000

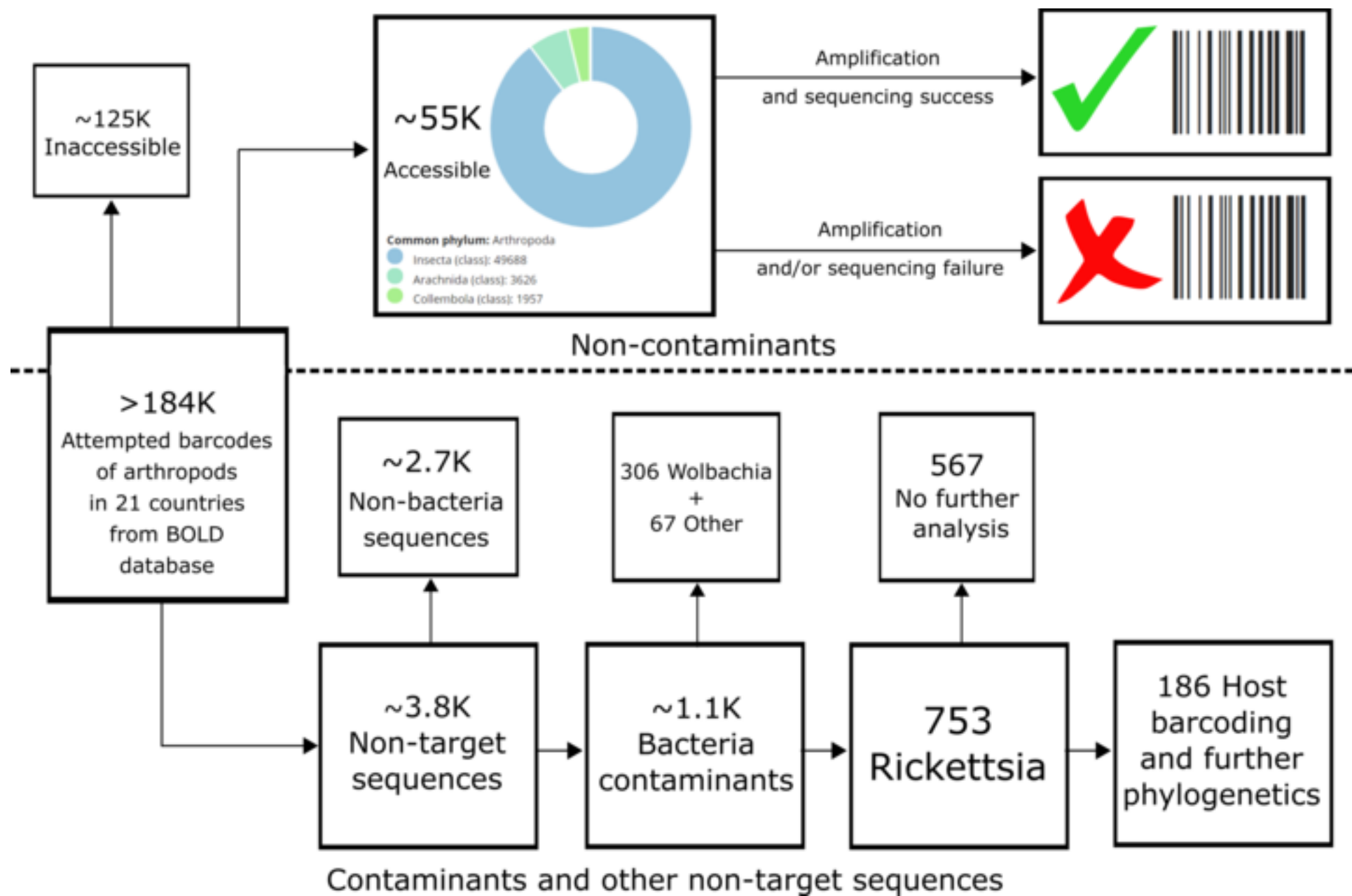
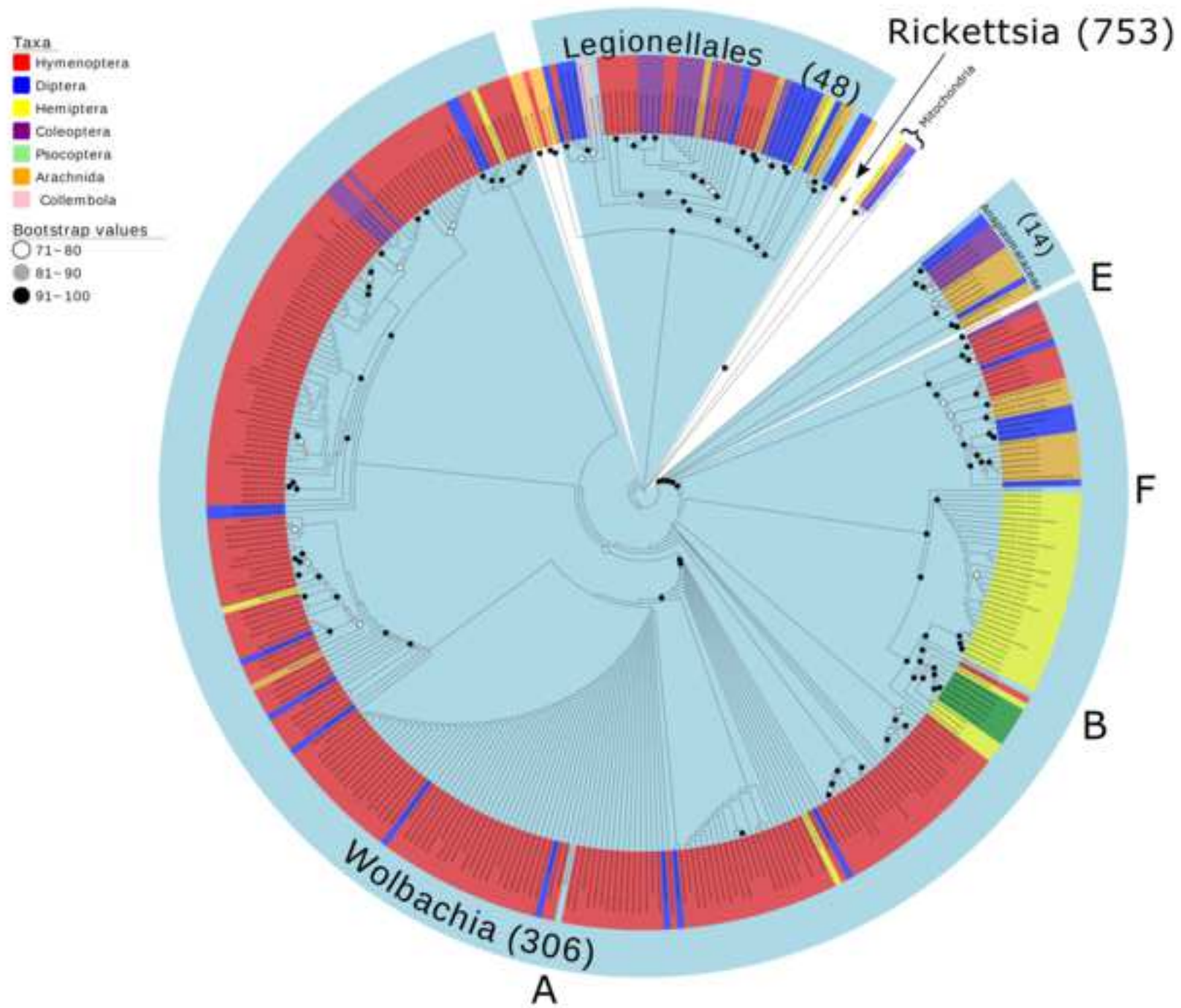


Figure 2

[Click here to access/download;Figure;Figure 2.png](#)





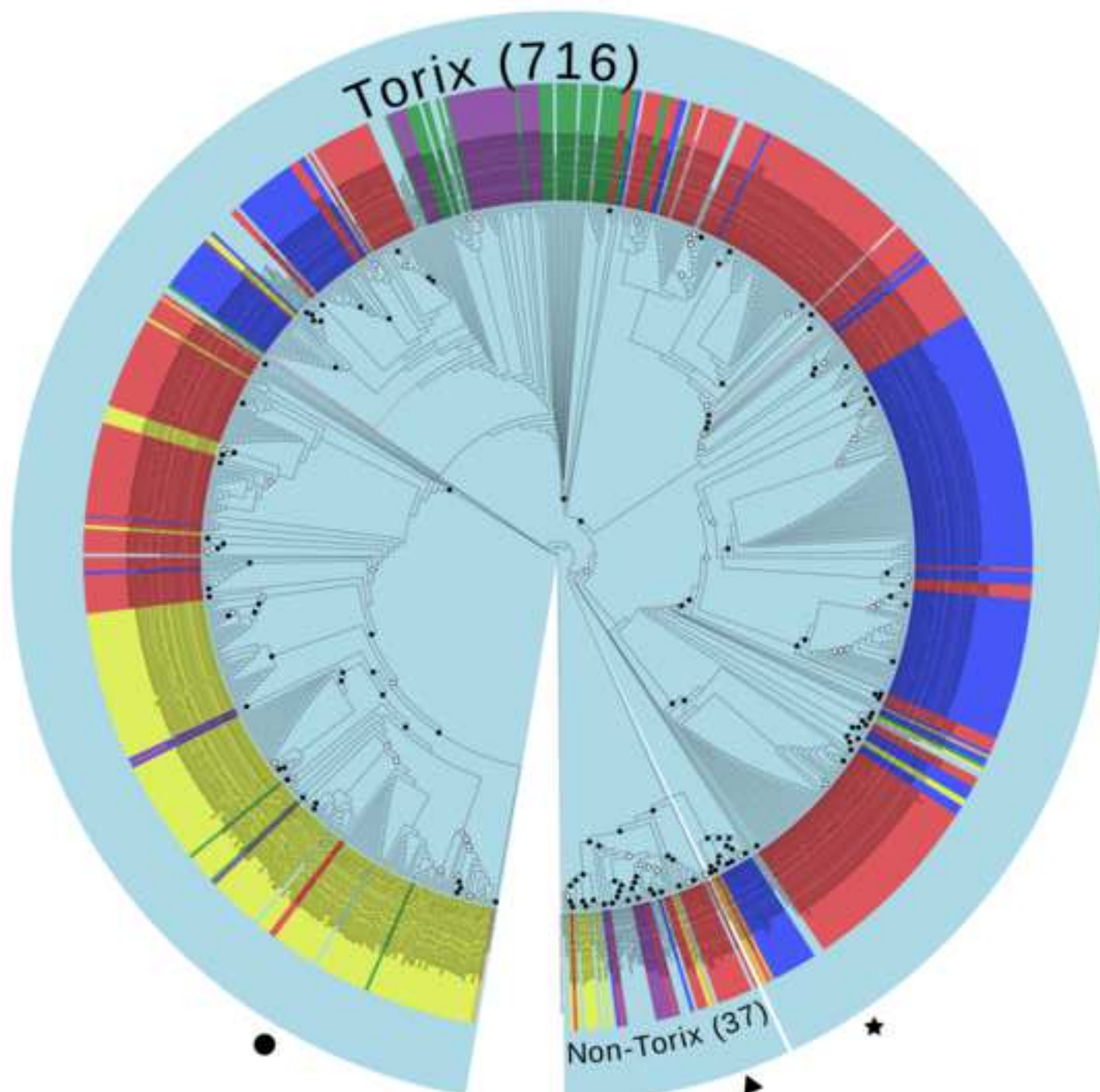
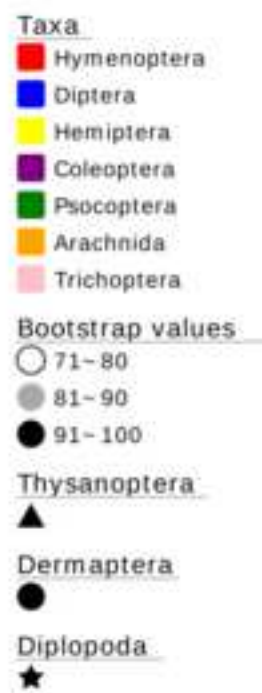




Figure 4

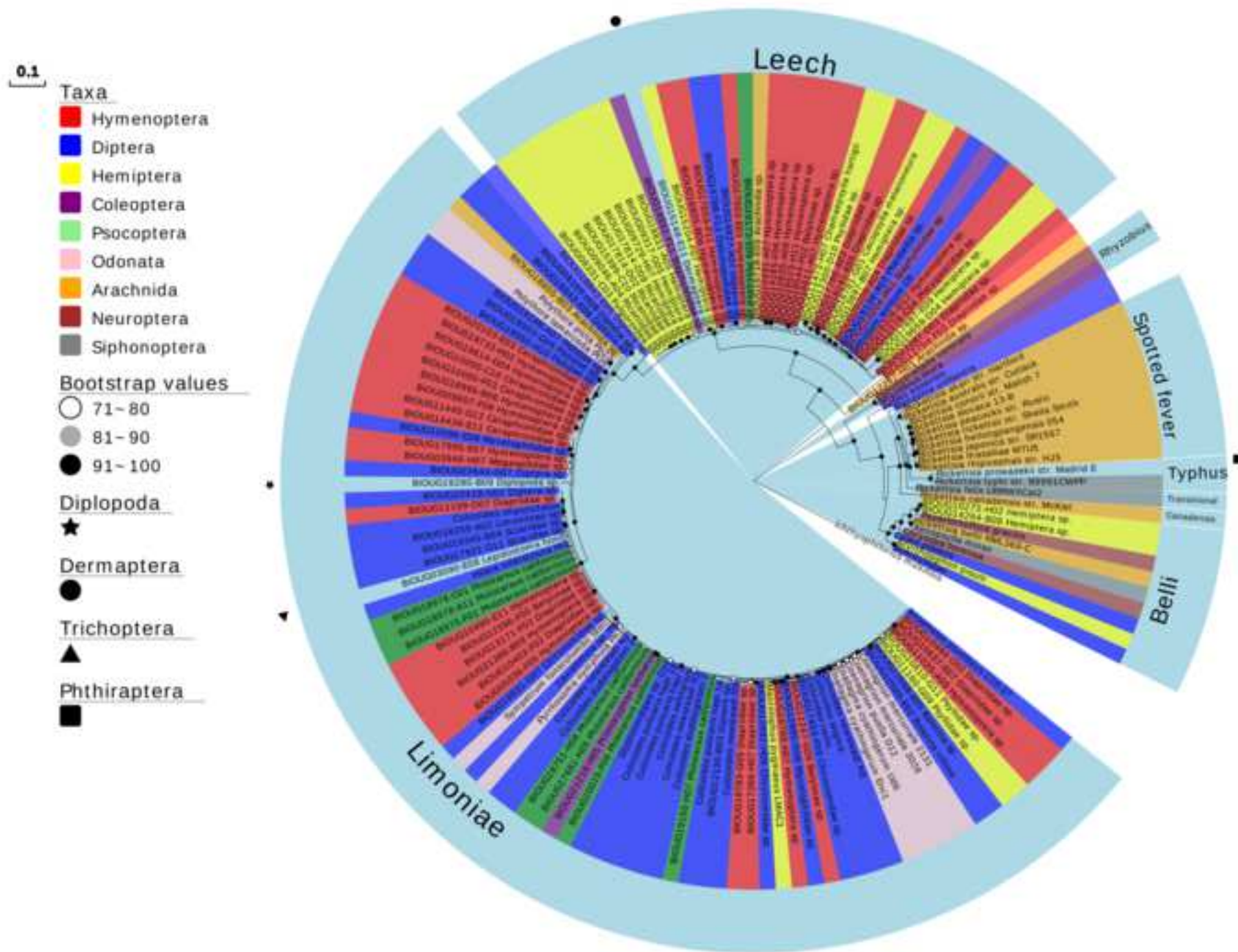
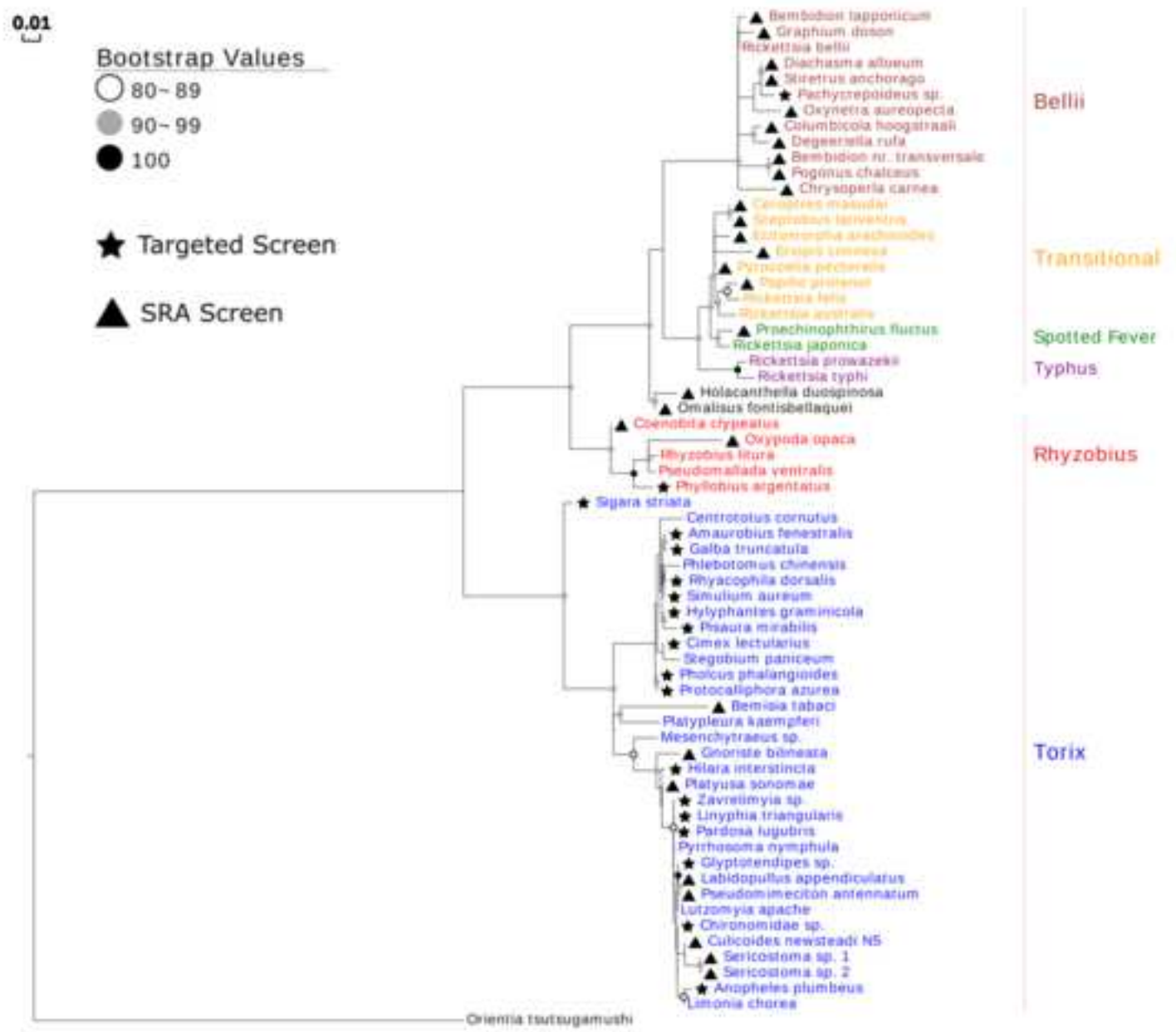
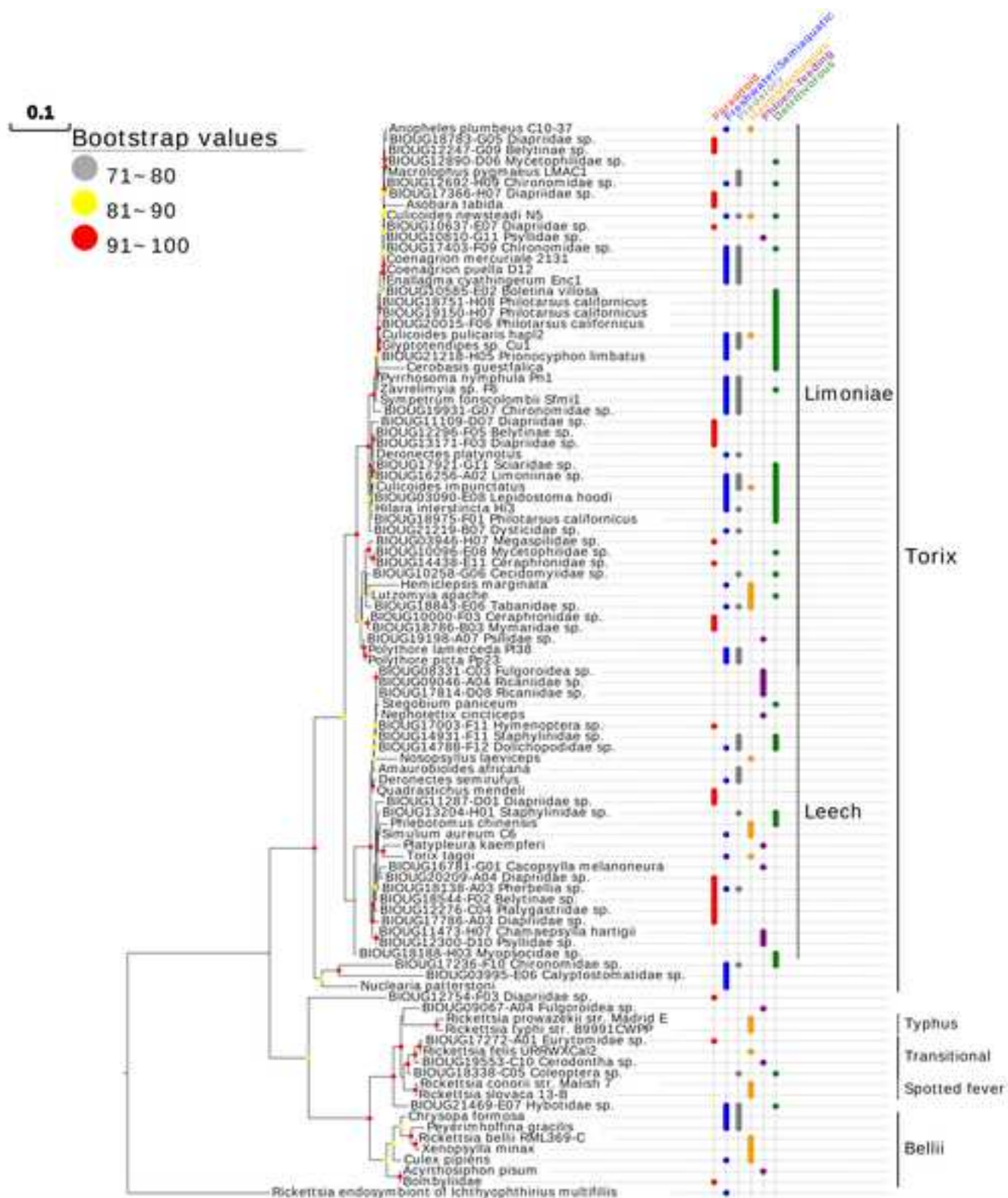
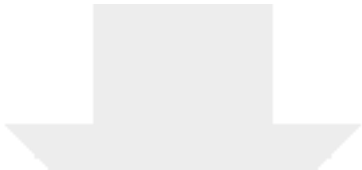


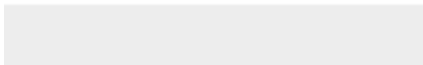
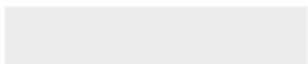
Figure 5

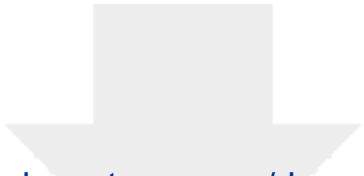




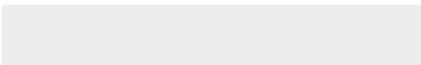



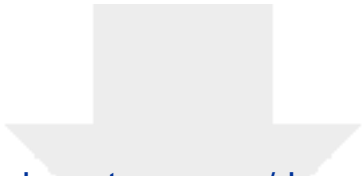
Click here to access/download  
**Supplementary Material**  
Additional file 1.docx






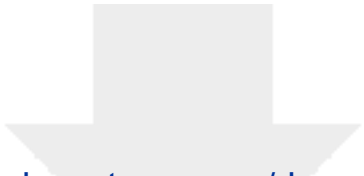
Click here to access/download  
**Supplementary Material**  
Additional file 2.7z






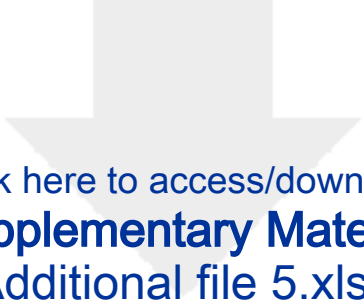
Click here to access/download  
**Supplementary Material**  
Additional file 3.docx



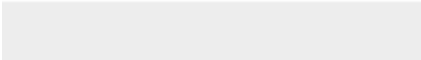



Click here to access/download  
**Supplementary Material**  
Additional file 4.docx

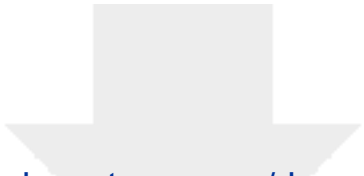




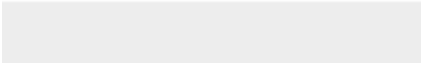

Click here to access/download  
**Supplementary Material**  
Additional file 5.xlsx

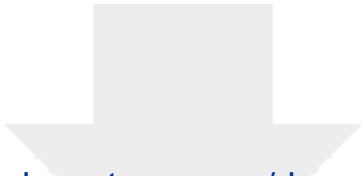







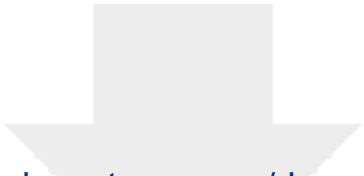
Click here to access/download  
**Supplementary Material**  
Additional file 6.docx



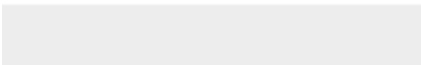



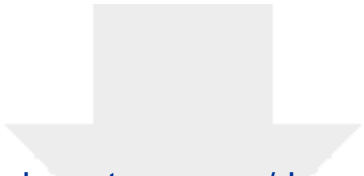
Click here to access/download  
**Supplementary Material**  
Additional file 7.docx






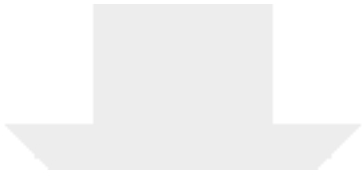
Click here to access/download  
**Supplementary Material**  
Additional file 8.docx






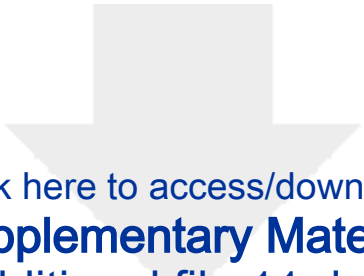
Click here to access/download  
**Supplementary Material**  
Additional file 9.pdf





Click here to access/download  
**Supplementary Material**  
Additional file 10.xlsx





Click here to access/download  
**Supplementary Material**  
Additional file 11.docx

